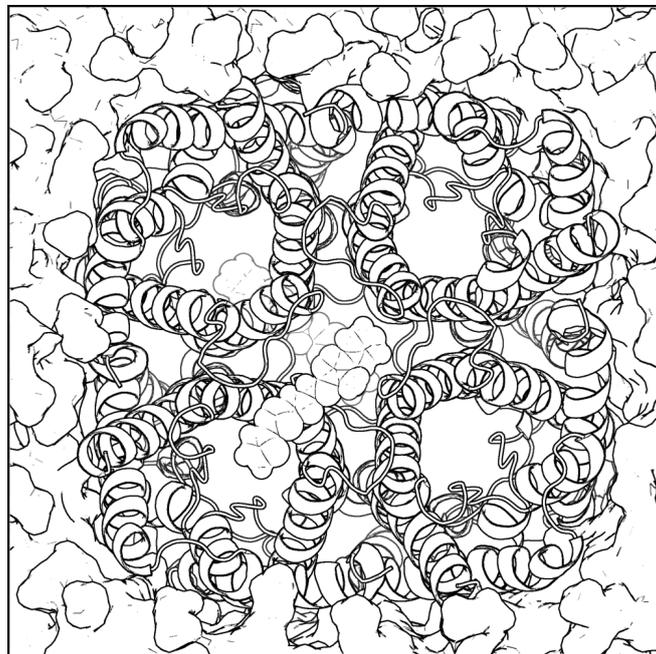


SÖREN WACKER

Computer-Aided Drug Design for Membrane Channel Proteins



Computer-Aided Drug Design for Membrane Channel Proteins

DISSERTATION

zur Erlangung des
mathematisch-naturwissenschaftlichen Doktorgrades
“Doctor rerum naturalium”
der Georg-August-Universität Göttingen

vorgelegt von

SÖREN WACKER
aus Stadthagen
Göttingen 2012

MITGLIEDER DES BETREUNGS-AUSSCHUSSES:

Prof. Dr. Bert L. de Groot (Begutachter)
Max Planck Institut für biophysikalische Chemie, Göttingen

Prof. Dr. Jörg Enderlein (Begutachter)
Georg-August-Universität Göttingen

Prof. Dr. Holger Stark
Max Planck Institut für biophysikalische Chemie, Göttingen

TAG DER MÜNDLICHEN PRÜFUNG: 07.08.2012

Computer-Aided Drug Design for Membrane Channel Proteins

DISSERTATION

for the award of the degree “Doctor rerum naturalium”
of the Georg-August-Universität Göttingen

submitted by

SÖREN WACKER
from Stadthagen
Göttingen 2012

MEMBERS OF THE THESIS COMMITTEE:

Prof. Dr. Bert L. de Groot (Reviewer)
Max Planck Institute for Biophysical Chemistry, Göttingen

Prof. Dr. Jörg Enderlein (Reviewer)
Georg-August-University Göttingen

Prof. Dr. Holger Stark
Max Planck Institute for Biophysical Chemistry, Göttingen

DATE OF THE ORAL EXAMINATION: 07.08.2012

This work is dedicated to
Sören from Nouna, Burkina Faso,
born in May 2012.

Part of this work have been published in the following articles:

Publications:

S. J. Wacker, W. Jurkowski, K. J. Simmons, C. W. G. Fishwick, A. P. Johnson, D. Madge, E. Lindahl, J.-F. Rolland, and B. L. de Groot. Identification of selective inhibitors of the potassium channel Kv1.1-1.2(3) by high-throughput virtual screening and automated patch clamp. *ChemMedChem*, Mar 2012.

S. Jelen, S. Wacker, C. Aponte-Santamaria, M. Skott, A. Rojek, U. Johanson, P. Kjellbom, S. Nielsen, B. L. de Groot, and M. Rützler. Aquaporin-9 protein is the primary route of hepatocyte glycerol uptake for glycerol gluconeogenesis in mice. *J Biol Chem*, 286(52):44319–25, 2011.

Hiermit bestätige ich, dass ich diese Arbeit selbstständig verfasst und keine anderen als die angegebenen Quellen und Hilfsmittel verwendet habe.

Göttingen, Juni 2012

Sören Wacker

Contents

1	Introduction	1
2	Theory and Concepts	5
2.1	Concepts in Drug Discovery	5
2.1.1	Sensitivity, Affinity, IC_{50} and Selectivity	6
2.1.2	Druglikeness - The Rule of Five	7
2.1.3	The Chemical Space	8
2.1.4	Ligand Similarity	8
2.1.5	The SMILES Notation	9
2.2	Introduction to Molecular Docking	9
2.2.1	Scoring Functions	11
2.2.2	Sampling Algorithms	13
2.2.3	Molecular Docking Programs	15
2.3	Evaluation of Virtual Screening Results	17
2.3.1	Enrichment	18
2.3.2	The ROC-curve	19
2.3.3	Metrics for Quality Assessment	20
2.3.4	Consensus Scoring	22
2.4	Thermodynamics of Ligand Binding	25
2.4.1	Thermodynamic Potentials	25
2.4.2	The Chemical Potential	26
2.4.3	The Binding Free Energy	27
2.4.4	Entropy-Enthalpy Decomposition	29
2.4.5	Potential of Mean Force	29
2.5	Molecular Dynamics Simulations	31
2.5.1	Approximations	32
2.5.2	General Simulation Conditions	34
2.5.3	Limitations	36
2.5.4	Calculation of the Binding Free Energy	38

3	Optimization of Molecular Docking	45
3.1	Introduction	45
3.2	Results	49
3.2.1	Detection of Possible Target Sites	49
3.2.2	Verification of Target Sites	50
3.2.3	Optimization of Molecular Docking	52
3.2.4	Prediction of Novel Active Compounds	53
3.2.5	Experimental Validation	54
3.2.6	Receptor Flexibility	57
3.3	Discussion	59
3.4	Summary	62
3.5	Methods	63
4	Inhibition of human Aquaporin 9	67
4.1	Introduction	67
4.2	Results	70
4.2.1	Homology Model of human Aquaporin 9	70
4.2.2	Single Pore Water Permeability Coefficients	70
4.2.3	High-Throughput Virtual Screening	74
4.2.4	Identification of Novel Inhibitors	77
4.2.5	Simulated Ligand Association	80
4.3	Discussion	82
4.4	Summary	84
4.5	Methods	86
5	Identification of First Active Compounds	89
5.1	Introduction	90
5.2	Results	93
5.2.1	Reproduction of the Crystal Structure	93
5.2.2	Structure Based Virtual Screening	94
5.2.3	Experimental Validation	95
5.2.4	Identification of First Inhibitors	99
5.3	Discussion	100
5.4	Methods	104
6	Outlook: Hit-Optimization based on Molecular Docking	105
6.1	Introduction	105
6.2	Implementation and Results	106
6.2.1	Compound Modification	106

6.2.2	The Docking Module	107
6.2.3	Final Scoring	108
6.2.4	Compound Efficiency	108
6.2.5	Quasi- <i>de novo</i> Design	110
6.3	Outlook	112
6.4	Methods	113
7	Conclusions	115
	Acknowledgements (German)	135
	Appendix	137

- ABC** ATP-binding cassette
ADMET adsorption, distribution, metabolic effects, excretion and toxicity
AQP Aquaporin
AROC area under the ROC-curve
AUAC area under the accumulation curve
ATP adenosine triphosphate
BAR Bennett's acceptance ratio method
BEDROC boltzmann enhanced discrimination of the ROC-curve
CADD computer aided drug design
CA classification accuracy
CDM chemically-defined medium
CNS central nervous system
CS consensus scoring
DFT density functional theory
eHiTS the electronic high-throughput screening environment
DtpT dipeptide and tripeptide transport protein
EF enrichment factor
EMA European Medicines Agency
FDA Food and Drug Administration
FEP free energy perturbation
G-protein guanine nucleotide-binding proteins
GAFF generalized amber force field
GPCR G-protein coupled receptor
HTS high-throughput screening
HTVS high-throughput virtual screening
IND Investigational New Drug
LBVS ligand based virtual screening
LIE linear interaction energy
MC monte carlo
MD molecular dynamics
MM molecular mechanics
MS multiple sclerosis
MWT molecular weight
NMR nuclear magnetic resonance
PMF potential of mean force
RIE robust initial enhancement
RMSD root mean square deviation
ROC receiver-operator-characteristic

SBD substrate-binding domain
SBVS structure based virtual screening
SF scoring function
SMILES simplified molecular-input line-entry system
TEA Tetraethylammonium
TI thermodynamic integration
TM transmembrane
VDW van der Waals
VS virtual screening
Vina Autodock-Vina
WHAM weighted histogram analysis method
wAROC weighted AROC
wAUAC weighted AUAC
US umbrella sampling

“Das Leben ist wert, gelebt zu werden, sagt die Kunst, die schönste Verführerin; das Leben ist wert, erkannt zu werden, sagt die Wissenschaft.”

Friedrich Nietzsche

1

Introduction

The ability of cells to sense and respond to chemical changes from each side of the plasma membrane is fundamental for all living organisms. Cells translate these changes, termed *signals*, into chemical changes. These so called *signal-transduction* processes are the basis for intra- and intercellular transfer of information and they enable cells to respond accordingly. Physiological signal-transduction can be achieved in various ways. In multicellular organisms, cells communicate by hormones and other chemical messengers, the concentration of ions and other solutes in the intra- or extracellular environment. On a molecular level, this means that mostly chemical changes activate or inactivate certain receptor proteins causing a signal cascade to be initiated, identified as the *response* to the particular chemical change. Respective important receptors are frequently located at the surface of the cell, namely the plasma membrane. More than 30 % of the human genome encode for integrated and associated membrane proteins. Their exposed position and role in the action and reaction of cells render membrane proteins important targets for therapeutical treatments. This is reflected by the fact that more than 50 % of all current drug targets are membrane proteins [96]. The predominant family addressed by current drugs are G-protein coupled receptors (GPCRs), a large family of eukaryotic transmembrane receptors that react to signals from the extracellular environment [12]. In response to extracellular changes, GPCRs generate intracellular responses mediated by

heterotrimeric guanine nucleotide-binding proteins (G-protein). Another class of integrated membrane proteins are membrane channels which facilitate the permeability of lipid membranes for certain solutes and water according to the chemical gradient. They are important for countless processes in the human body including neural conduction, the cardiac action potential and osmotic balance. Altogether, membrane channels are ideal drug targets with an enormous potential for future pharmaceutical treatment [96]. However, they are conceptually different from other drug targets as enzymes or GPCRs. Evolutionary, the pores of membrane channel proteins are not optimized for high-affinity ligand binding. This might be the reason why only five per cent of all current drug targets are membrane channels [12].

Structural biology enables the direct search and the modelling of functional modifiers which complement the structure and the chemistry of a particular receptor [22]. A widely applied method in the field of drug discovery is *molecular docking* that relies on the three dimensional receptor and ligand structures in order to predict the structure of receptor-ligand complexes [66]. A common application of molecular docking in the field of drug design is the screening of virtual compound databases for the identification of putatively active compounds (actives), compounds that activate or inactivate a biological target. Hereby, the affinity of the ligand, usually in terms of the standard binding free energy, is estimated by evaluating a so called scoring function (SF). This function quantifies the interactions of receptor and ligand in the complex. The generation of scoring functions generally includes fitting to experimental data including a broad scope of receptor-ligand complexes to yield generality. Unfortunately, in the set of available liganded protein structures membrane channels are in the minority. Therefore they are rarely used for the training of scoring functions, rendering molecular docking against membrane channels challenging.

In principle, molecular docking tries to assess the standard binding free energy of complex formation. For the estimation of receptor-ligand interactions the free energy is the most important thermodynamic characteristic [43, 44, 140]. In general, it describes the driving forces of basically all biological process as e.g. the folding of proteins, osmotic forces and in particular the formation of receptor-ligand complexes. Knowing the basic physics involved in these processes theoretically enables to calculate the corresponding binding free energy. However, the complexity of biological systems renders the exact calculation impossible for typical biological systems. The equations that describe such systems are analytically and computationally untractable. Despite that for many systems it is possible to construct a discretized, virtual *model system* that reflects the relevant,

inherent properties of the real *target system* [136]. Modern computer based techniques use parameter based model systems and several computational and mathematical “tricks” for the generation of structural ensembles corresponding to time series of structures representing a dynamical processes. This time series is in general referred to as *simulation*. Such simulations are expected to correspond to the dynamics of the target system. When such a biological process can be modeled and covered by the simulated timescales, simulations can be used to approximate the binding free energy involved in these processes. Nevertheless, the size of typical biologically systems and the timescales on which e.g. drug binding takes place (timescales up to milliseconds) cause an enormous amount of computational power rendering an application on thousands or millions of compounds impossible.

In this work, I combine computational methods as molecular docking and all-atom molecular dynamics simulations to explore the inhibition of membrane channel proteins by non-covalent association of small chemical compounds. The first chapter is focused on the optimization of molecular docking techniques targeting the chimeric potassium channel $K_V1.1-(1.2)_3$ and shows how contemporary molecular docking algorithms can be optimized for the efficient prediction of potassium channel inhibitors. In the second chapter, I explore the inhibition of the human water channel protein hAQP9 by combining various computational methods as molecular docking and all-atom molecular dynamics simulations. This study revealed the location of the interaction site of inhibitors. Furthermore, the complete binding process of a known inhibitor was simulated. Chapters 5 and 6 describe the *status quo* of ongoing studies with perspectives for future engagements. Chapter 5 covers the initial phase in a drug discovery endeavor starting with a crystal structure and no knowledge about small molecule inhibitors. The latter study is focused on the phase after the successful identification of active compounds and covers the process of compound optimization.

In summary, membrane channel proteins hold an enormous potential for the development of novel pharmaceutical treatments. Furthermore, selective inhibitors of individual membrane channels are valuable for the study of the physiological role of membrane channels including their involvement in (human) diseases. In addition, little is known about the actual binding process of ligands to membrane channel proteins in general.

2

Theory and Concepts

2.1 Concepts in Drug Discovery

Drug discovery can be defined as the process in which chemical compounds with activity against a target or a function are identified. Desired effects could be the suppression of gene products, inhibition of an enzymatic reaction, the interference with a signaling cascade, inactivation of transport proteins or the blocking of channel proteins. The initial identification of active compounds usually requires a reliable functional assay and a collection of compounds for screening. Then, compounds that show sufficient activity in this initial screen (*hits*) are evaluated on the basis of potency, specificity, toxicity and efficacy in animal models and other properties to select *lead* compounds [127], which will enter the clinical phase. The phase between hit identification and lead selections is called the hit-to-lead phase. Currently applied hit-identification strategies range from knowledge-based approaches, which use literature-derived molecular entities, endogenous ligands or biostructural information to quasi 'brute-force' methods such as combinatorial chemistry or high-throughput screening (HTS). The dominant and the most widely applicable technique for the identification of lead compounds is HTS [12, 115], an experimental screening technique based on robotics where large numbers of different compounds are screened in a time as short as possible and at reasonable costs. Per day, 1,000 – 100,000 individual assays can

be carried out in a typical HTS setup [75, 76, 126]. Usually, 50,000 – 1,000,000 compounds are tested in one single screen. The results obtained in HTS depend significantly on the type of assay used in the screen. Sills et al. [116] showed that different types of active compounds are identified by different assay types.

2.1.1 Sensitivity, Affinity, IC_{50} and Selectivity

The ability of an entity like a cell to respond to an external signal is called the *sensitivity*. The higher its sensitivity is, the lower is the threshold of the signal to cause a response. Hereby, the sensitivity can be increased by cooperative effects or –in the case of receptor-ligand interactions– high affinities of the ligands to the receptors. The affinity is quantified as the association constant K_a or its reciprocal counterpart the dissociation constant K_d . For a receptor-ligand complex reaction



where R is the receptor, L the ligand and RL the complex, K_a and K_d are determined by the equilibrium concentrations of the receptor C_R^0 , the ligand C_L^0 and the complex C_{RL}^0 or by the on- and off-rates k_{on} and k_{off} :

$$K_a = \frac{k_{on}}{k_{off}} = \frac{C_{RL}^0}{C_R^0 C_L^0} = \frac{1}{K_d} \quad (2.2)$$

The effectiveness of a molecule that inhibits a certain biological target, function or reaction, can be measured quantitatively by the half-maximal inhibitory concentration, the IC_{50} . Regarding the binding of molecules to its receptors the IC_{50} is the ligand concentration where the concentration of liganded and unliganded receptors is equal. Often, the IC_{50} is converted to the pIC_{50} :

$$pIC_{50} = -\log_{10}(IC_{50}) \quad (2.3)$$

The IC_{50} is not a direct indicator of the binding affinity. However, for competitive agonists (inhibitor) and antagonists (substrate) both can be related by the Chen-Prusoff equation:

$$K_a = IC_{50} \left(1 + \frac{C_S}{C_{S,50}} \right)^{-1} \quad (2.4)$$

where C_S is the concentration of the substrate and $C_{S,50}$ the substrate concentration where the activity of the receptor is half-maximal when no inhibitor is present. The *selectivity* of a ligand for a certain receptor measures how specific a ligand binds a certain receptor with respect to other receptors or causes a certain response. The selectivity for a certain receptor with respect to other receptors

can be quantified by the fraction of the binding affinities. Sometimes, the level of inhibition of different receptors or phenotypes at a fixed ligand concentration is used to estimate affinity and/or the specificity.

2.1.2 Druglikeness - The Rule of Five

As the number of compounds in libraries of large pharmaceutical companies used in HTS was approaching 1 million, logistic obstacles and cost issues made this library size an upper limit for most companies [12]. After the realization that the quality for reliable and information-rich biological readouts cannot be obtained using ultra-high synthesis techniques, many research organizations subsequently scaled back their large scale production rates and focused on smaller but structurally diverse compound libraries. The content of present compound libraries in pharmaceutical companies is more driven by the question of what is useful than what is possible. Accordingly, the outcome of early combinatorial chemistry approaches has been widely replaced by smaller contents that are structurally focused to compounds which are considered to be *drug-like* or *lead-like*, meaning molecules that structurally resemble marketed drugs or lead compounds. In 2001, Lipinski *et al.* [81] set a landmark for the estimation of the oral applicability of compounds by the definition of *the rule of 5*. A set of properties that nowadays, has widely been taken as the definition of *drug-likeness*. Based on a distribution of calculated properties among several thousand drugs, the rule of 5 predicts poor adsorption or permeation properties when there are more than 10 H-bond acceptors, more than 5 H-bond donors, a molecular weight (MWT) of more than 500 Dalton and a calculated LogP of more than 5. Lead-like compounds, in contrast, have a lower MWT (around 300 Dalton) and have fewer H-bond donors and acceptors. Notably, Lipinski suggested that compound classes that are substrates for biological transporters are exceptions to the rule, because these compounds are transported actively across membranes. Therefore, the general structural constraints that are necessary to in order arrive at its target receptor e.g. diffuse through the lipid-bilayers, are not required for these class of compounds. Also antibiotics, antifungals, vitamins and cardiac glycosides are exceptions to the rule of 5 [81]. Therefore, the accordance with the rule of 5 is not a guarantee for good metabolic properties and an exception not an absolute exclusion criterion. In any case, the rule of 5 concentrates research at a property space with reasonable possibility of oral activity and thus makes labor-intensive studies of drug metabolisms more efficient.

2.1.3 The Chemical Space

The growing number of different chemical entities in the databases that are used in the drug development process raise the question about the relative number of these compounds and how they compare to each other. The set of all possible chemical compounds is frequently conceptualized as the *chemical space* or the *chemical universe*, in analogy to the cosmic universe, and can be defined as the set of all possible molecular structures. It is widely accepted that the chemical space is huge, but the estimation of the absolute number of its elements varies by several orders of magnitude. Bohacek *et al.* [13] estimated the number of compounds with a maximum number of 30 carbon, nitrogen, oxygen and sulfur atoms to exceed 10^{60} , whereas Ertl [35] considered the number of organic molecules that can be synthesised with currently know methods and estimated it to be between 10^{20} and 10^{24} . An extensive review about the different estimations of the size of the chemical space was published by Medina-Franco *et al.* in 2008 [89]. However, for medicinal chemistry a much smaller fraction of compounds will be relevant, since the majority of these structures will reveal a poor pharmacokinetic profile, i.e. poor adsorption, distribution, metabolic effects, excretion and toxicity (ADMET) properties.

2.1.4 Ligand Similarity

Comparing molecules is a challenging task. A widely applied concept in chemical informatics are chemical *fingerprints*. The fingerprint of a molecule is a sequence of bits or boolean array that is generated with respect to structural features of the molecule. The assessment of the ligand's similarity then breaks down to the comparison of bitstrings, assuming that the similarity of the bitstrings contains information about the similarity of the underlying molecular structures. The similarity of the fingerprints can then be assessed by applying the Tanimoto metric [104], also called the Tanimoto coefficient, distance or similarity. The Tanimoto similarity $T(a, b)$ of two bit sequences is defined by

$$T(a, b) = \frac{N_c}{N_a + N_b - N_c} \quad (2.5)$$

where N_a and N_b are the total numbers of bits of each string and N_c the number of bits that is present in both strings, referred to as the intersection of a and b . When there is no overlap between a and b , $T(a, b)$ becomes zero. When a and b are identical $T(a, b)$ becomes one. Two molecules are considered similar, when the corresponding Tanimoto coefficient of the molecules fingerprints is larger than 0.7.

2.1.5 The SMILES Notation

The simplified molecular-input line-entry system (SMILES) is a chemical structure specification that uses one dimensional ASCII strings to encode chemical structures. Originally invented by Arthur and David Weininger in the 1980s, it was further modified mainly by Daylight Chemical Information Systems Inc. Typically, multiple valid SMILES-strings can be written for a molecule. For example, CCO, OCC and C(O)C all specify the structure of ethanol. Atoms are encoded by the standard chemical abbreviation in square brackets. For a subset of organic molecules (N, O, P, S, F, Cl, Br, and I) the square brackets can be omitted. Hydrogen atoms can explicitly be added, otherwise the canonical number of hydrogen atoms is assumed. A specific protonation state can be provided by adding an H, the number of hydrogen atoms, a number of +/- for atomic charges, e.g.: [NH4+] for a ammonium ion and [Co+3] or [Co+++] for a cobalt 3+ ion. Between aliphatic atoms single bonds are assumed unless other bond types are specified. “=” stands for a double bond and “#” for a triple bond. Aromaticity is represented by lower case letters. The connectivity in ring systems is encoded by digits, e.g. “c1ccccc1” for benzene. For systems with more than 9 rings the “%” character has to be put before the ring label. Branches are represented by parentheses e.g.: “C(C)(C)(C)C” for 2,2-dimethylpropane. Configuration around double bonds is specified using the characters “/” and “\”: “F/C=C/F” for the trans- and “F/C=C\F” for the cis-configuration. The stereochemistry of molecules with stereo centers can be specified by “@”, for example L-alanine can be written as “N[C@@H](C)C(=O)O” and D-alanine “N[C@H](C)C(=O)O”. The specifier “@@” indicates that, when viewed from nitrogen along the bond to the chiral center, the sequence of substituents hydrogen (H), methyl (C) and carboxylate (C(=O)O) appear clockwise.

2.2 Introduction to Molecular Docking

A key method for the prediction of the structures of receptor-ligand complexes in the lead and drug discovery process is *molecular docking* [93]. This technique was first applied in late 1980s and is widely used as a virtual screening tool in the early stage of the drug development process. Furthermore, it has been invaluable for the understanding of receptor-ligand interactions. In the following, I will highlight important aspects of molecular docking with respect to the work presented here. The docking process involves three phases. The first phase, the *sampling*, covers the generation of ligand configurations and orientations of a

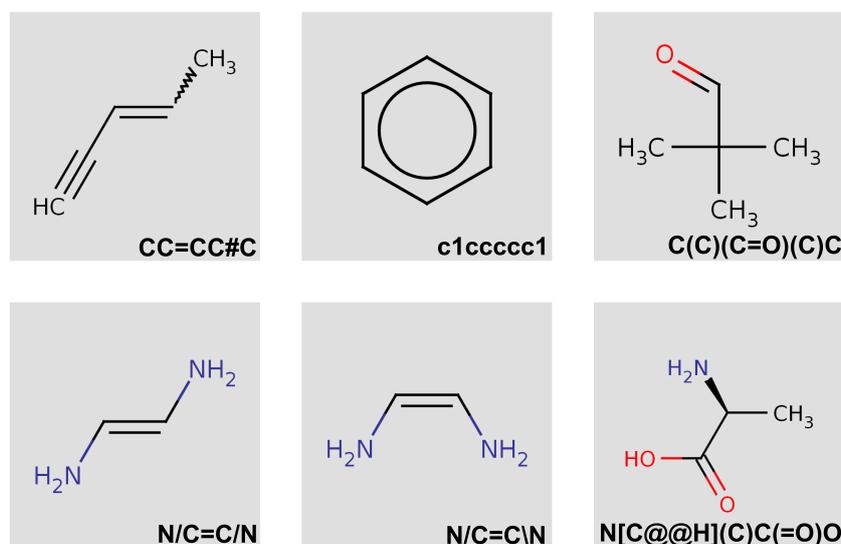


Figure 2.1: Six different chemical structures and corresponding SMILES-strings.

ligand relative to the target binding site. These are referred to as *poses*. When receptor flexibility is taken into account, the sampling also involves the variation of the receptor configuration. In the second phase, the *scoring*, a so called *docking score* is calculated as an estimate of the ligand binding affinity or activity (Section 2.1.1). When docking is applied to screen virtual compound libraries, the compounds are ranked according to the best scored poses. This process is called *ranking*. The score is calculated by evaluating the *scoring function*, that often represents the binding free energy of the complex. Hereby, the complexity of the receptor ligand interaction is immensely reduced. Most of the contemporary scoring algorithms are focused on enthalpic terms, whereas molecular associations are also driven by entropic effects. Often docking programs used simplified structural representations and reduce if not neglect protein flexibility as well as the participation of solvent molecules in binding. Additionally, most docking programs assume a certain static protonation state and consider a fixed distribution of charges among the atoms. The lengths and, except for the torsions of *rotatable bonds*, angles between covalently bonded atoms are kept fixed [121]. However, the benefit of molecular docking has been demonstrated in many studies. In the following, the most important components of molecular docking will be covered in more detail.

2.2.1 Scoring Functions

The *scoring function* is one of the central concepts in molecular docking. This function enables a docking algorithm to rapidly describe and quantify the interactions between ligand and receptor. During the sampling phase the docking algorithm produces different ligand configurations and orientations within the target site and assigns a score by evaluating the scoring function. Hereby, an ideal scoring function would provide the lowest scores for the energetically most favorable receptor ligand configurations. Assuming that these configurations represent the interactions that mainly promote the ligand binding, they give direct insight into the underlying molecular mechanisms. An excellent overview over a broad spectrum of scoring functions is given in [94]. There are mainly three different types of scoring functions used:

Force-field based scoring functions are designed based on underlying physical interactions such as van der Waals (VDW) interactions, electrostatic interactions as well as bond stretching, bending and torsional interactions. The force field parameters are usually derived by both fitting to empirical data and *ab initio* calculations. A typical force-field based scoring function is implemented in the DOCK algorithm whose energy function is the sum of VDW and coulombic energy contributions:

$$E = \sum_{i,j \neq i} \left(\frac{A_{ij}}{r_{ij}^{12}} - \frac{B_{ij}}{r_{ij}^6} + \frac{q_i q_j}{\epsilon r_{ij}} \right) \quad (2.6)$$

where A_{ij} and B_{ij} are VDW parameters, q_i and q_j the charges and r_{ij} the distance between the particles i and j . ϵ is the dielectric constant [90]. Equation 2.6 does not include the energetic costs of desolvation which is a many body interaction term and depends on the chemical environment. In order to account for the desolvation, further terms are usually added based on the solvent-accessible surface area of the ligand and possibly the receptor [58]. When energy terms of VDW and Coulomb interactions are used in a scoring function, they need to be significantly empirically weighted, in part, to account for the difference between energies and free energies [19, 43, 121], and in part, to account for the different methods used to calculate the different terms.

Empirical scoring functions estimate the binding free energy ΔG of a receptor-ligand complex by a sum of weighted energy terms:

$$\Delta G = \sum_i w_i \Delta G_i \quad (2.7)$$

The energy terms G_i can represent VDW and electrostatic interactions, hydrogen bonding strength, entropy changes, hydrophobic interactions or desolvation

energies and other contributions. The weights w_i are derived by fitting to known experimental data of a training set. In 1994 Böhm developed an empirical scoring function consisting of hydrogen bonds, polar interactions, the lipophilic contact area between ligand and receptor and the number of rotatable bonds in the ligand. The weights were calibrated with a dataset of 45 protein-ligand complexes [14]. This scoring function was further improved when Eldridge *et al.* [33] developed the ChemScore scoring function that includes terms for hydrogen bonds, metal atoms, lipophilic contacts as well as the number of rotatable bonds in the ligand. With the number of different terms in an empirical scoring function it becomes more and more difficult to avoid the double counting of specific interactions. The applicability of empirical scoring functions may depend on the data used in the training set. Empirical scoring functions that are fitted to larger training sets promise to be more generally applicable.

Knowledge based scoring functions use terms that weight the receptor-ligand complexes by the occurrence frequencies of particle-particle pairs in a database of known complexes [119]. The idea behind the knowledge based scoring function is as follows: Supposing large numbers of different particles (e.g. amino acids or atoms) were somehow to distribute themselves in a gas phase at temperature T if the interactions are purely pairwise, the distributions can be described by the equilibrium pairwise density $\rho_{ij}(r)$ between any two particle types $i, j = 1, 2, \dots$ at distance r . In this case, the interaction free energy, $w_{ij}(r)$, can be calculated from the observed densities by the inverse Boltzmann relation:

$$w(r) = -k_B T \ln \left(\frac{\rho_{ij}(r)}{\rho_{ij,0}} \right) \quad (2.8)$$

where $\rho_{ij}(r)$ is the pair density of a particle pair at distance r and $\rho_{ij,0}$ the pair density of a reference state where the interatomic interactions are zero [58, 119]. Since these potentials are extracted from the structures rather than from attempting to reproduce known binding affinities by fitting, and because the training structural database can be large and diverse, the knowledge-based scoring functions are quite robust and relatively insensitive to the training set [58]. Because of the pairwise interaction scheme the knowledge based scoring functions can be as fast as the empirical scoring functions. However, atoms in protein-ligand complexes are not particles in the gas phase and the pair frequencies are not independent from each other. Therefore, the calculation of accurate reference states $\rho_{ij,0}$ is a challenging task in the development of knowledge based scoring functions.

Hybrid scoring functions are implementations of mixtures of the different flavors of scoring functions. They combine for example force field terms and empirical

energy terms. This is done, for example, in the program eHiTS [146, 147] which is described in more detail in section 3.5. Notably, all currently applied scoring functions require a significant degree of empirical fitting. Therefore, scoring functions are not necessarily generally applicable to all kinds of drug targets and should be benchmarked and possibly optimized against special or rare forms of receptors.

2.2.2 Sampling Algorithms

Molecular docking algorithms can be classified by their search algorithms, which are applied to predict the complex structure. Search algorithms split up in *global optimization search algorithms* that aim to sample systematically the complete search space and *guided progression search algorithms* that focus their search to promising parts of the search space. When treating ligand (and receptor) flexibility, the global searches suffer from a combinatorial explosion, because even for very small compounds with few rotatable bonds N_r the number of possible conformation N_C is extremely large:

$$N_C = \prod_{i=1}^{N_r} \prod_{j=1}^k \frac{360}{\Phi_{i,j}} \quad (2.9)$$

Here k is the number of increments and $\Phi_{i,j}$ the size of the increments. Therefore, most contemporary sampling techniques use guided progression searches, that reject unfavorable conformations and therefore greatly reduce the number of conformations [72].

Considering the treatment of flexibility of the ligand/receptor there are two groups of search algorithms. The *rigid body search algorithms* do not take flexibility into account. They basically solve a 6 dimensional (3 translations, 3 rotations) two body optimization problem. Because of the low dimensionality they can work extremely fast. Rigid body docking algorithms usually rely on *fast shape matching algorithms* that take into account the geometrical overlap of receptor and ligand. The shape matching can be the only criterion for the calculation of the docking score or combined with interaction and desolvation terms as it is done in ZDOCK [100]. Fast shape matching algorithms are also used by *flexible docking algorithms*, the second class of docking algorithms, which take either ligand or receptor flexibility or both into account and thereby require more computational power [30].

With respect to the degree of flexibility, the class of flexible docking algorithms splits up in several subclasses. There are docking algorithms that take into account only the ligand flexibility and treat the receptor as rigid. Other algorithms

also take receptor flexibility into account starting with alternative side chain rotamers of the receptor amino acids and ending with the full flexibility of the receptor.

In principle, the docking problem can be addressed by applying energy minimization techniques that perform local optimizations, for example the *steepest decent* algorithm. The problem with these techniques is that they do not explore the configuration space exhaustively and the results depends highly on the initial placement of the ligand[101]. *Simulated annealing* aims to avoid getting trapped in local minima. This algorithm uses a stochastic optimization procedure that allows phase space transitions contrary to the local energetic gradient with a certain probability. Hereby, the acceptance rates are determined by a *metropolis criterion* and a successively decreasing temperature parameter T is applied. The acceptance probability for such an algorithm can be expressed as

$$P_{\text{accept}} = \min(1, \exp(-\frac{\Delta E}{T})) \quad (2.10)$$

where ΔE is the change in potential energy and T a free parameter with the same unit as ΔE . However, differences in results obtained using simulated annealing on different starting position show that similar problems as with local minimization techniques occur in practice. Another technique to avoid the trapping in local minima is the systematic variation of the input parameters as done in *genetic algorithms*. They perform several runs for the search and therefore are computationally expensive. A systematic guided progression search with low redundancy is the *incremental fragmentation* as used in FlexX. It is regularly applied to target rigid receptors. Here the ligand is cut into fragments which are successively docked into the target site. Thereby, the ligand is reconstructed. Usually, the first fragment serves as starting point or “anchor” for the reconstruction. Therefore, these algorithms are also referred to as *anchor-and-grow* methods[30]. However, there are also different subtypes of this technique, that pursue different strategies for both the fragmentation and the reconstruction, as done in the program eHiTS (Section 3.5).

There are also exotic methods available such as the search in Fourier space or the *distance geometry* method. The latter affords representations of the conformational space of the ligand in form of a matrix that contains constraints for all atom-atom distances. This matrix is a complete but highly redundant description of the conformational space of the ligand. Furthermore, only a small subset of distance matrices which are consistent with the constraints represent meaningful conformations. The translation from the distance space to euclidian coordinates is computationally expensive, but euclidian representations are required for the

calculation of the interaction energy with the receptor. Therefore this method is computationally very expensive [101]. Since in guided progression searches the scoring function is sometimes evaluated during the sampling procedure, a sharp separation of sampling and scoring is not always possible.

2.2.3 Molecular Docking Programs

In the following, two molecular docking approaches are presented in detail. I focus on the programs Autodock-Vina (Vina) and FlexX (or LeadIT). These programs are conceptually different and reflect the variety of molecular docking approaches.

Autodock-Vina (Vina)

AutoDock Vina [121] (version 1.0.2) – hereafter termed Vina – is an open source docking suite, using an iterative local search algorithm and several runs starting from random conformations. For the local search, a quasi-Newton method is used that does not calculate the Hessian matrix of the potential surface explicitly. A succession of steps is performed that consist of mutations and local optimizations. Each step is accepted according to a Metropolis criterion (Section 2.2.2). The arguments of the function that is optimized are the location and orientation of the ligand as well as the torsion angles for all rotatable bonds. Several runs starting from random initial arguments are performed. The number of the runs is varied with respect to the apparent complexity [121]. Significant minima are then combined and used for structure refinement and clustering. The general form of the scoring function used in Vina is

$$c = \sum_{i < j} f_{t_i, t_j}(r_{ij}) \quad (2.11)$$

where summation goes over all pairs of atoms that can move relative to each other. These are interactions between atoms that are separated by three covalent bonds. Each atom is assigned a type t_i . The interaction function $f_{t_i, t_j}(r_{ij})$ between atoms at interatomic distance r_{ij} is symmetric. Hereby, f_{t_i, t_j} is actually a function of the atomic surface distance $d_{ij} = r_{ij} - R_{t_i} - R_{t_j}$ of the radii R : $f_{t_i, t_j}(r_{ij}) =$

$h_{t_i,t_j}(d_{ij})$. The following energy terms are used in the scoring function of Vina:

$$\begin{aligned}
 E_{\text{Gauss}_1}(d) &= \exp\left(-\left(\frac{d}{0.5 \text{ \AA}}\right)^2\right) \\
 E_{\text{Gauss}_2}(d) &= \exp\left(-\left(\frac{d-3 \text{ \AA}}{2 \text{ \AA}}\right)^2\right) \\
 E_{\text{Repulsion}} &= \begin{cases} d^2, & \text{if } d < 0 \\ 0, & \text{if } d \leq 0 \end{cases} \\
 E_{\text{hydrophobic}} &= \begin{cases} 1, & \text{if } d < 0.5 \text{ \AA} \\ \text{linear interpolation} & \\ 0, & \text{if } d > 1.5 \text{ \AA} \end{cases} \\
 E_{\text{H-bond}} &= \begin{cases} 1, & \text{if } d < -0.7 \text{ \AA} \\ \text{linear interpolation} & \\ 0, & \text{if } d > 0 \text{ \AA} \end{cases}
 \end{aligned}$$

The hydrophobic and the hydrogen-bond term have a pairwise linear form. All interactions are cut off at a distance $r_{ij} = 8 \text{ \AA}$. c can be expressed as sum of intra- and intermolecular interactions. Then the predicted binding free energy G is calculated from the intermolecular interactions c_{inter} by

$$G(c_{\text{inter}}) = \frac{c_{\text{inter}}}{1 + w_r N_r} \quad (2.12)$$

Each energy term is associated with a weight. The energy function is not evaluated every time the ligand adopts a new pose. Instead, Vina calculates a grid map for each atom type from the fixed part of the receptor.

FlexX/LeadIT

FlexX [101] uses an incremental approach for the flexible docking of ligands. At first the algorithm selects a connected and rigid part of the ligand as the *base*. The base is chosen automatically and placed into the defined target site. Next, the ligands are incrementally reconstructed. During this reconstruction process, fragments of the ligand are successively fit to the base fragment in all possible conformations. The best of these placements (by the scoring function) are used for the next reconstruction step. The scoring function implemented in FlexX has

the form:

$$\begin{aligned}
 \Delta G &= \Delta G_0 + \Delta G_{\text{rot}} N_{\text{rot}} \\
 &+ \Delta G_{\text{coul}} \sum_{\text{coul}} f(\Delta R, \Delta \alpha) \\
 &+ \Delta G_{\text{hbond}} \sum_{\text{hbond}} f(\Delta R, \Delta \alpha) \\
 &+ \Delta G_{\text{arom}} \sum_{\text{arom}} f(\Delta R, \Delta \alpha) \\
 &+ \Delta G_{\text{lipo}} \sum_{\text{lipo}} f'(\Delta R)
 \end{aligned}$$

The argument $\Delta R = R - R_i - R_j - 0.6 \text{ \AA}$ is the distance of the atoms minus the radii of the individual atoms R_i and R_j and an additional offset of 0.6 \AA . The terms ΔG_i correspond to ideal geometries. The functions $f(\Delta R, \Delta \alpha)$ penalize deviations from these geometries. During the reconstruction procedure the scoring function is evaluated for the selection of the “best” solutions. Hereby, the different interaction terms are weighted differently. Optionally, another set of weights can be used for a final evaluation of the reconstructed ligands. Then, the final scores are used to rank the set of docked ligands.

2.3 Evaluation of Virtual Screening Results

In this section the different evaluation methods of molecular docking or virtual screening experiments are explained. Emphasis was placed on the analysis of the enrichment of known active compounds in a subset of top scored compounds.

One essential measure for the performance of a molecular docking algorithm is the reproduction of native binding modes defined by a threshold in the root mean square deviation (RMSD). Although, the rating quality of the RMSD is problematic for small and large molecules, it has been widely used as criterion for the definition of success or failure of docking algorithms [58].

A second criterion for the performance of a molecular docking algorithm is its ability to predict the binding affinity of different ligands. Because the scale of docking scores is not always in the range of experimental data, often the correlation between docking scores x_i and experimental data y_i in form of the Pearson correlation coefficient C_P is considered.

$$C_P = \frac{\sum_{k=1}^N (x_k - \langle x \rangle) (\langle y_k - \langle y \rangle)}{\sqrt{(\sum_{k=1}^N (x_k - \langle x \rangle)^2) (\sum_{k=1}^N (y_k - \langle y \rangle)^2)}} \quad (2.13)$$

where N is the number of ligands or complexes, and x_i and y_i the corresponding scores and experimental values. C_P is useful to measure a linear correlation,

but the correlation between scores and experimental values is not linear. In that case, it is better to project both the scores and experimental values to ranks, and calculate the correlation between the ranks. This is exactly what the Spearman correlation coefficient C_S stands for.

2.3.1 Enrichment

In structure based virtual screening (SBVS) molecular docking is used to screen databases of compounds in order to identify active compounds. For this purpose, the docking scores has to separate active compounds from inactive compounds. In order to test a particular docking approach it is possible to dock a library of compounds with known active and inactive compounds to the known binding site. Then the success of the docking can be estimated by the *enrichment*, the fraction of the active compounds in a subset of top scored compounds. Another approach for the benchmark of docking algorithms is the similarity of the generated ligand poses with the poses found in an experimentally derived structure. When the score represents the binding free energy, active compounds should be scored lower than non binders. The tendency to score active compounds differently leads to a shift in the relative probability distributions (Figure 2.2). A screen of the top ranked compounds should then find preferentially actives. One can estimate the quality of a molecular docking algorithm by preparing a test library with active and inactives and monitor the number of active compounds with respect to the docking score. This can be done in different ways as described in the following.

The *enrichment* $\epsilon(x_0)$ is defined as the accumulated rate of active compounds within the top x_0 per cent of a ranked list that contains both known active compounds, and inactive compounds or decoys. It is bounded at the points (0,0) and (1,1) –or (100,100) when interpreted as percentages. The most direct way to plot the enrichment is realized by the *accumulation curve* $\epsilon(x)$ with $x = \{0, 1\}$. Herein, it is referred to the *enrichment plot* when $\epsilon(x)$ is plotted on a logarithmic scale. The *enrichment factor* $\xi(x_0)$ is defined as the fraction of the enrichment of a ranked list and the expected enrichment of a randomly sorted list at a certain point x_0 :

$$\xi(x) = \frac{\epsilon(x)}{f_a x} \quad (2.14)$$

where $f_a = \frac{N_a}{N}$ is the fraction of active compounds N_a in a total of N compounds. The upper limit of $\xi(x)$ depends on the absolute number of compounds and the fraction of active compounds f_a . The run of both curves depends strongly on the threshold that defines active and inactive compounds as illustrated in figure 2.3.

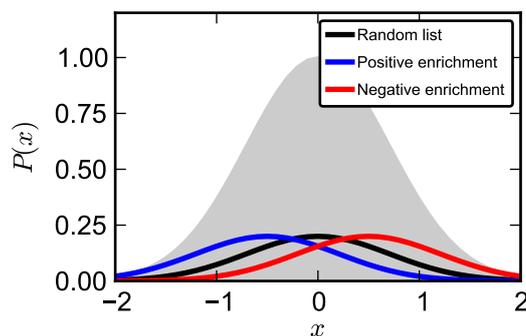


Figure 2.2: Schematical illustration of the separation of active and inactive compounds with respect to the docking scores. The grey area refers to an assumed distribution of random decoys. The black line corresponds to a scoring function that does not shift active compounds with respect to the decoys and therefore does not lead to any enrichment of positively identified active compounds. The blue curve corresponds to a shift towards lower scores and the red curve towards higher scores.

The enrichment factor is important for the estimation of the optimal fraction of a docked and ranked library of compounds with unknown activity to be screened experimentally.

2.3.2 The ROC-curve

Another way to monitor a virtual screening (VS) result is the receiver-operator-characteristic (ROC) curve, which is widely applied in other fields [38]. In general, ROC-curves are parameter curves that monitor the true-positive (tp) rate on the Y-axis and the false-positive (fp) rate on the X-axis. These rates depend on a discrete classifier. For each value of the classifier a pair (tp, fp) is generated corresponding to a single point in the ROC space. Here the classifier is the threshold score/rank of a list of scored/ranked compounds. The rates tp and fp correspond to the fraction of identified actives and identified inactives. The points (0,1) and (1,0) correspond to identifiers that perfectly classify the compounds. Hereby (0,1) means all actives are identified as actives and all inactives are identified as inactives. (1,0) means that all actives are identified as inactives and vice versa. In figure 2.3, the curves of $\epsilon(x)$, $F_\epsilon(x)$ and the ROC-curve are shown with respect to different thresholds that define active and inactive compounds. Whereas $\epsilon(x)$ and $F_\epsilon(x)$ vary significantly, the ROC-curves are relatively robust. As indicated in figure 2.3 the ROC-curve of a perfect ranking

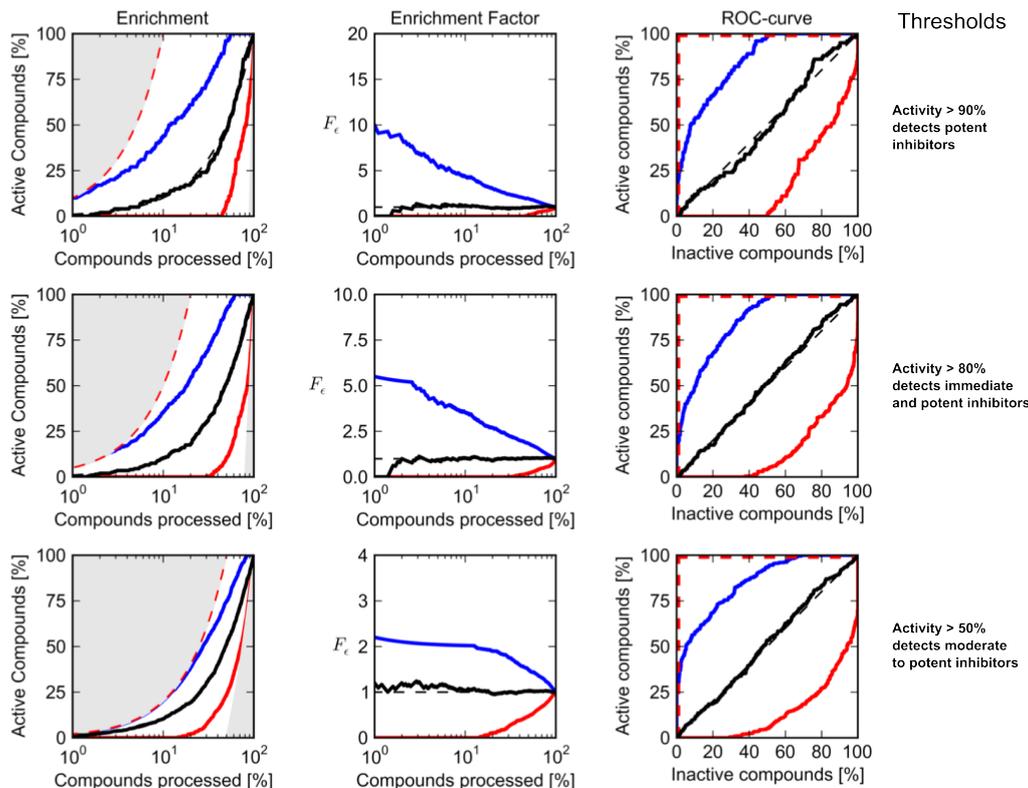


Figure 2.3: $\epsilon(x)$, $F_\epsilon(x)$ and the ROC-curve are shown with respect to different thresholds that defines actives and inactives. The curves correspond to distributions that are shown in figure 2.2. The same color coding has been used. Each set of black, blue and red curves corresponds to the same list of scored compounds. The dashed lines correspond to a perfect separation (red) and an ideal random distribution (black).

(red dashed lines) looks always the same. This property suits the ROC-curve to serve as a robust estimator for the quality of molecular docking algorithms.

2.3.3 Metrics for Quality Assessment

There are several metrics available for the measure of performance of VS results. Truchon *et al.* [60] published a detailed discussion of contemporary metrics. One frequently used metric is the area under the accumulation curve (AUAC).

$$\text{AUAC} = \int_0^1 dx \epsilon(x) \quad (2.15)$$

The same can be done with the ROC-curve. The area under the ROC-curve (AROC) has been frequently used to measure VS performance [60] and is widely

applied in other fields. The AROC is relatively independent on the ratio of active compounds. It is bounded by 0 and 1. Hereby, 1 corresponds to a perfect positive enrichment, where all actives are scored lower than all the inactives. An ideal random ranking would correspond to 0.5 and an absolute negative enrichment where all inactives are scored lower than all the actives corresponds to 0. It can be interpreted as a probability that an active compound is ranked before inactive compound or decoy.

$$\text{AROC} = \int_0^1 dx \text{ROC}(x) \quad (2.16)$$

A problem with the AUAC and AROC is that these metrics do not distinguish between early and late recognition of active compounds as illustrated in figure 2.4. A hypothetical VS result could rank half of the actives very low and the other half very high. Then, both AUAC and AROC would give a value of 0.5 as in the case of a complete random ranking, although there is a meaningful difference between these situations. In order to overcome this limitation it is possible to weight the different contributions of the area under the curves with respect to the argument x by a weighting function $w(x)$.

$$\text{wAUAC} = \frac{\int_0^1 dx w(x) \text{AUAC}(x)}{\int_0^1 dx w(x)} \quad (2.17)$$

$$\text{wAROC} = \frac{\int_0^1 dx w(x) \text{ROC}(x)}{\int_0^1 dx w(x)} \quad (2.18)$$

The weighting function can be an exponential function $w(x) = \exp(-\alpha x)$. This form has the advantage that the extend of weighting can be controlled by a single parameter α . Due to the weighting the wAUAC is not necessarily bounded by 0 and 1, causing that a perfect enrichment is associated with an arbitrary number which depends on the weighting function. A useful modification of the wAUAC is the so called boltzmann enhanced discrimination of the ROC-curve (BEDROC) metric as introduced by Truchon *et al.* [60]:

$$\text{BEDROC} = \frac{\text{wAUAC} - \min(\text{wAUAC})}{\max(\text{AUAC}) - \min(\text{AUAC})} \quad (2.19)$$

The BEDROC metric discriminates between early and late recognition of true positives and is bound to 0 and 1. Figure 2.4 show how the values of the different metrics depend on the threshold for active compounds. The BEDROC therefore is particularly suitable to assess a scoring method's ability to identify true actives in a small selection of top ranked compounds.

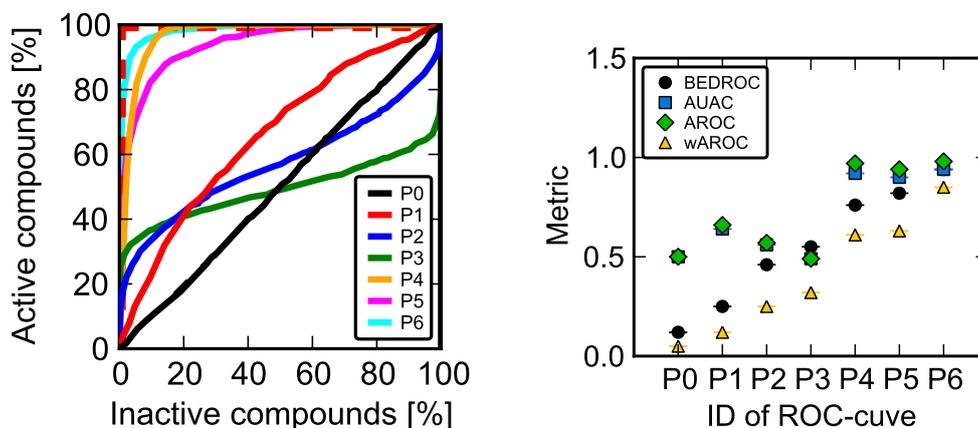


Figure 2.4: Different hypothetical shaped ROC-curves (left) and the corresponding values of different the metrics (right).

2.3.4 Consensus Scoring

A possible method for the enhancement of virtual screening (VS) results is consensus scoring (CS) first applied in molecular docking by Charifson *et al.* [20]. The main idea of CS is to combine different VS results in order to obtain better agreement with the experimental results or higher hit-rates. Nowadays, consensus scoring (CS) is widely applied for the enrichment of virtual libraries, the prediction of binding poses or binding affinities. In my work, I used CS solely for the enhancement of the enrichment of virtual screenings. The enhancement of the enrichment using CS has been demonstrated in several publications [20, 97, 137]. A comprehensive review about CS was published by Feher in 2006 [39].

The bandwidth of strategies and techniques used for CS is broad. With *strategy*, I mean the way how different VS approaches are applied. For example, a compound can be docked with two different programs or the pose of a ligand generated in one program can be evaluated with the scoring function of another program (*rescoring*). Instead of combining the values of SFs also the individual terms of SF can be combined [117]. In contrast, with *technique*, I mean the mathematical way of combining different scores (averaging, minimum, maximum, weighting, etc.). Which strategy is followed and the technique which is applied depends on the goal of the study (identifying the correct binding pose, maximize the enrichment, find the most affine compound, etc.) and the specific conditions (number of compounds to be evaluated, available computer power, desired number of compounds in final set, etc.). Finally, with *approach*, I mean

the combination of a particular strategy with a certain technique.

The first consensus approach used in a VS study was intersection based [20]. It involved scoring compounds with multiple scoring functions and taking the intersection of the top $N\%$ scored compounds. Only compounds which occur in the top $N\%$ of all applied scoring functions are selected. A feature of this technique is that the intersection of top ranked compounds by definition is smaller than the original list of compounds. Therefore, the number of compounds in the final list becomes smaller the more SFs are applied. A modification of the intersection technique overcomes this (sometimes) disadvantageous property by assigning “votes” to the compounds. If a candidate is predicted to be on the top $N\%$ by a certain SF, then it gets a “vote” from that scoring function. The final score of a candidate compound is the number of votes gathered from all the scoring functions, which may range from 0 to the total number of scoring functions. This approach is widely known as *rank-by-vote*.

Another technique of CS is to build a linear combination of the individual SF values. When this linear combination is simply the average of all scores I refer to it as *rank-by-num*. Alternatively, the compounds can be ranked by the minimal (*rank-by-min*) or maximal (*rank-by-max*) scores of the scoring functions [135] what would be a special case of a weighting technique. These techniques are useful when the scoring functions assign comparable numbers to the compounds, e.g. when all scores reflect the absolute or relative binding free energy of the system. Sometimes, the scores are on very different scales or even reflect different entities, for example the potential of mean force (PMF), and the binding free energy or an arbitrary number without physical meaning. In this case the consensus can be build according to the ranks of the compounds. When the average of the ranks is used for the consensus, I refer to it as *rank-by-rank* [135]. The rank-by-rank technique may also be interpreted as a compromise between the intersection based technique and rank-by-number. It is also possible to combine different schemes in a VS approach. For example, it may be useful to screen a compound library against different receptor structures using different scoring functions and then apply the rank-to-min technique with respect to different receptor structures first and than apply the rank-to-max technique in a second step.

In order to explain why CS works, Wang and Wang [135] performed an idealized computer experiment with a hypothetical set of 5000 compounds, and analysed the relationship between the *hit-rates*, the rate of correctly identified actives, and the number of SFs used for the consensus. They assumed that the value of the SF is the activity of the ligand plus a random number and observed that the number

of false positives and false negatives decrease with increasing number of SFs. The enrichment in the top 100 scored ligands also increased continuously with the number of SFs, when using the rank-by-number or rank-by-rank technique. When using the rank-by-vote, they observed a steadily decreasing number of hits with increasing number of applied SFs, as stated by Carifson *et al.* [20]. Finally, Wang and Wang [135] concluded that:

“[...] the consensus scoring outperforms any single scoring [function] for a simple statistical reason: the mean value of repeated samplings tends to be closer to the true value.”

Apart from this general mathematical reason there may be also structural reasons involved that originate from the structural knowledge about receptor and the compounds. At least for ligand based scoring techniques, it was observed that scoring functions tend to provide more similar rankings for active compounds than for inactives [4]. It is widely accepted that the most benefit from CS can be obtained when the individual contributions (i.e. the factors which are combined) perform well on the particular target and when the individual factors are not correlated. The involvement of factors without predictive power, in general, decreases the informative value of the consensus. Whereas the use of correlated factors may lead to an overestimation of certain contributions e.g. the hydrophobic contacts. Applying a certain SF to other docking programs can lead to inaccuracies and errors, because distances between ligand and receptor atoms can vary when using different docking programs and the applied SF can be sensitive to these differences. Therefore, the consensus of results from different docking experiments and the consensus of several scoring functions applied for a single docking experiment (rescoring) are conceptually different [39].

Z-scores

When the scales of the individual factors (SFs values or SF terms) are too different for the construction of a reasonable consensus score, it is possible to convert these factors into z-scores and to build the consensus with the corresponding z-scores [87].

$$z_i = \frac{x_i - \mu}{\sigma} \quad (2.20)$$

where μ is the mean value and σ the standard deviation of a population x_i of N values. Using an arbitrary input the z-scores project it to a distribution with a mean value of zero and a standard deviation of one. Therefore, the z-scores are of similar order of magnitude and can be used for a CS. In this work, I have always used z-scores when applying a CS technique.

2.4 Thermodynamics of Ligand Binding

This section provides the thermodynamic framework of ligand binding following the works of Gilson and Zhou [44, 140]. Central concepts mentioned in section 2.1.1 as e.g. the IC_{50} and the standard binding free energy ΔG_0^b are derived from the basic laws of thermodynamic. The unfeasibility of the exact calculation of ΔG_0^b is rationalized and prospects for its approximation are derived, namely in terms of the potential of mean force (PMF).

2.4.1 Thermodynamic Potentials

The equilibrium state of a thermodynamic system is fully described by the fundamental thermodynamic relation

$$dU = TdS - pdV + \sum_i \mu_i dN_i \quad (2.21)$$

which describes the change of the internal energy U of a system as a function of the entropy S , the volume V and the number of particles N_i . Here the temperature T , the pressure p and the chemical potential μ_i of particle species i are used. U can be denoted as:

$$U = TS - pV + \sum_i \mu_i N_i \quad (2.22)$$

The partial derivatives of $U(S, V, N)$ according to S, V and N are:

$$\left(\frac{\partial U}{\partial S}\right)_{V, N_i} = T \quad (2.23)$$

$$\left(\frac{\partial U}{\partial V}\right)_{S, N_i} = -p \quad (2.24)$$

$$\left(\frac{\partial U}{\partial N_i}\right)_{V, S} = \mu_i \quad (2.25)$$

what demonstrates that the system is fully determined when $U(S, V, N)$ is known. Therefore, U is called a *thermodynamic potential* with the *natural variables* S, V and N . In experiments, it is usually easier to control the temperature than the entropy. Therefore, it is more convenient to describe the system as a function of T instead of S . The differential of the free energy $F = U - TS$ is

$$dF = -SdT - pdV + \sum_i \mu_i dN_i \quad (2.26)$$

Under conditions of constant temperature T , volume V and number of particles N , $F(T, V, N)$ becomes a thermodynamic potential, that is minimized

at the thermodynamic equilibrium. At experimental conditions of constant temperature and constant pressure the Gibbs free energy (the free enthalpy) $G(T, p, N) = U + pV - TS = F - TS$ becomes a thermodynamic potential that is minimized at the thermodynamic equilibrium. The differential of $G(T, p, N)$ is:

$$dG = -SdT + Vdp + \sum_i \mu_i N_i \quad (2.27)$$

2.4.2 The Chemical Potential

The thermodynamic potentials $U(S, V, N)$, $F(T, V, N)$ and $G(T, p, N)$ provide measures of the stability of a system at a thermal equilibrium: the lower they are, the higher the stability. In a microscopic system, high stability corresponds to a high probability of occupancy, what corresponds to a high statistical weight Q , that is known as the partition function. The free energy F and the partition function are related by

$$F = -k_B T \ln(Q) = -k_B T \ln(Q), \quad (2.28)$$

where $\beta = (k_B T)^{-1}$ is the inverse of the product of the absolute temperature T and the Boltzmann constant $k_B = 1.3806 \cdot 10^{-23} \text{ JK}^{-1}$. When a system with the energy function $E(\vec{x})$ can be described by a set of canonical coordinates \vec{x} , the partition function Q of the system is given by the integral

$$Q = \xi \int d\vec{x} e^{-\beta E(\vec{x})} \quad (2.29)$$

Here the constant ξ is inserted to render Q unitless. It is convenient to integrate over the external translations, that lead to a factor V , and rotations, that lead a factor $8\pi^2$. Then eq. 2.29 becomes an integral over the internal coordinates \vec{x}' :

$$Q = 8\pi^2 V \xi \int d\vec{x}' e^{-\beta E(\vec{x}')} \quad (2.30)$$

In the following, I will omit the prime for the internal coordinates.

When the microstates separate into two non-overlapping macrostates A and B , the corresponding partition functions Q_A and Q_B are of the form of eq. 2.30 with the integral going over the corresponding microstates. The probability of occupancy p_i for the state i is then proportional to the corresponding partition function Q_i :

$$p_i = \frac{Q_i}{Q} \quad (2.31)$$

which is the reason why the partition function can be interpreted as a statistical weight. Considering now a system with N_i non-interacting and identical particles

i with single particle partition function Q_i , the partition function of the total system becomes:

$$Q = \frac{Q_i^{N_i}}{N_i!} \quad (2.32)$$

The factor $N_i!$ accounts for the indistinguishability of the N_i particles. The expression for the chemical potential in section 2.1.1 can now be derived by inserting eq. 2.32 in eq. 2.28 and applying Stirling's approximation $\ln N! \approx N \ln(N) - N$:

$$F_i = -k_B T \ln(Q) \approx -k_B T N_i \ln \frac{Q_i e}{N_i}, \quad (2.33)$$

where e is Euler's number. Now, the chemical potential of the particles i can be derived by the partial derivative of F with respect to N_i :

$$\mu_i = \left(\frac{\partial F}{\partial N_i} \right)_{T, V_0, N_{i' \neq i}} = -k_B T \ln \left(\frac{Q_i}{N_i} \right) = -k_B T \ln \left(\frac{Q_i/V}{C_i} \right) \quad (2.34)$$

where V is the volume and C_i the concentration of molecules i . The total chemical potential of the solvated system becomes then:

$$\mu = -k_B T \sum_i \ln \left(\frac{Q_i/V}{C_i} \right), \quad (2.35)$$

a quantity that is independent on the volume because each Q_i contains the term V (compare eq. 2.30). With this general expression of the chemical potential a relation between the free energy of binding and the association constant K_a can be derived as presented in the next section.

2.4.3 The Binding Free Energy

Considering the reaction from section 2.1.1 between the ligand L and the receptor R :



The concentrations of unbound ligand is $[L]$, the concentration of the free receptors is $[R]$ and the concentration of the complex is $[RL]$. The equilibrium constant K_a of these reaction is defined as:

$$K_a = \frac{C_{RL}^0}{C_R^0 C_L^0} \quad (2.37)$$

where C_R^0 and C_L^0 are the equilibrium concentrations of the unbound receptor and the unbound ligand, C_{RL}^0 the equilibrium concentration of the bound complex. Upon binding under constant pressure and constant temperature, the Gibbs free

energy ΔG is the difference between the chemical potentials of the complex μ_{RL} and the sum of chemical potentials of the receptor μ_R and the ligand μ_L :

$$\Delta G = \mu_{RL} - \mu_R - \mu_L \quad (2.38)$$

Inserting the expression for the chemical potential from eq. 2.35 results in:

$$\Delta G = -k_B T \ln \left(\frac{C_R C_L}{C_{RL}} \frac{(Q_{RL}/V)}{(Q_R/V)(Q_L/V)} \right) \quad (2.39)$$

When all species R , L , and RL are present at a certain standard concentration C_0 (usually 1M) the relation for the standard binding free energy is obtained as

$$\Delta G_0^b = -k_B T \ln \left(C_0 \frac{(Q_{RL}/V)}{(Q_R/V)(Q_L/V)} \right) \quad (2.40)$$

It is worth noting that the value of ΔG_0^b depends on the choice of the standard concentration. When converting a given K_a to a standard binding free energy ΔG it is important to adjust the standard concentration accordingly [44, 140]. Recognizing that $\Delta G = 0$ for the equilibrium establishes the identification of K_a by

$$K_a = \frac{C_{RL}^0}{C_R^0 C_L^0} = \frac{(Q_{RL}/V)}{(Q_R/V)(Q_L/V)} \quad (2.41)$$

multiplication with the standard concentration C_0 and using eq. 2.40 results in:

$$\Delta G_0^b = -k_B T \ln (C_0 K_a) \quad (2.42)$$

and therefore

$$K_a = e^{-\beta \Delta G_0^b} \quad (2.43)$$

Remembering that K_a is the ratio of the on-rate k_{on} and the off-rate k_{off} of the binding reaction of eq. 2.1, the physical meaning of ΔG_0^b becomes obvious:

$$\frac{k_{on}}{k_{off}} = e^{-\beta \Delta G_0^b} \quad (2.44)$$

When $\Delta G_0^b = 0$ then the on-rate and the off-rate are equal. The concentration will stay constant and the system is in equilibrium with respect to the concentrations. When $\Delta G_0^b < 0$, k_{on} is larger than k_{off} and the association of the ligands prevails the dissociation. Therefore, the concentration of the complex C_{RL} increases. When $\Delta G_0^b > 0$, k_{on} is smaller than k_{off} and the dissociation of the ligands prevails the association.

2.4.4 Entropy-Enthalpy Decomposition

The binding free energy can be decomposed in an enthalpic term that describes the interaction strenght of the receptor and the ligand, and an entropic term that describes the relative degree of uncertainty of the bound and the unbound state. Knowing ΔG_b^b , the entropy component of binding can be derived by

$$\Delta S_b = - \left(\frac{\partial \Delta G_b}{\partial T} \right)_{p, N_i} \quad (2.45)$$

With eq. 2.42 we get

$$S = k_B \ln \left(\frac{C_R C_L K_a}{C_{RL}} \right) + k_B T \left(\frac{\partial \ln(K_a)}{\partial T} \right)_p + k_B T \left(\frac{\partial}{\partial T} \ln \left(\frac{C_R C_L}{C_{RL}} \right) \right) \quad (2.46)$$

When the experiment is performed at constant pressure p and constant number of particles N_i , the concentrations C_i may change due to changes in the volume V and have to be considered as a function of V . The last term in eq. 2.45 can be reduced to a single factor

$$k_B T \left(\frac{\partial}{\partial T} \ln \left(\frac{C_R C_L}{C_{RL}} \right) \right) = k_B T \left(\frac{\partial}{\partial T} \ln \left(\frac{1}{C_V} \right) \right) = \frac{k_B T}{V} \left(\frac{\partial V}{\partial T} \right)_{p, N_i} = -k_B T \kappa, \quad (2.47)$$

where κ is the thermal expansion coefficient of the solution. Pure water has a $\kappa \approx 2.6 \times 10^{-4} \text{ K}^{-1}$ at 25°C , what would correspond to a value of $-0.08 k_B$ [140]. Therefore, the last term in eq. 2.45 can usually be neglected. Replacing the concentrations by the standard concentration C_0 yields the standard binding entropy:

$$\Delta S_b^0 = k_B \ln(C_0 K_a) + k_B T \ln \left(\frac{\partial \ln(K_a)}{\partial T} \right)_{p, N_i} \quad (2.48)$$

The binding enthalpy can be obtained as

$$\Delta H_b^0 = \Delta G_b^0 + T \Delta S_b^0 = k_B T^2 \left(\frac{\partial \ln(K_a)}{\partial T} \right)_{p, N_i}, \quad (2.49)$$

which has no dependence on the standard concentration C_0 . Therefore, a change of the standard concentration shifts the absolute values of ΔG_b^b and ΔS_b^0 , but not of ΔH_b^0 .

2.4.5 Potential of Mean Force

For a system consisting of solvent and solutes with continuous external “solvent” coordinates \vec{X} and internal “solute” coordinates \vec{x} and energy function $E(\vec{X}, \vec{x})$, the partition function eq. 2.29 can be written as

$$Q = \xi \int d\vec{X} d\vec{x} J(\vec{X}, \vec{x}) e^{-\beta E(\vec{X}, \vec{x})} \quad (2.50)$$

The Jacobian factor $J(\vec{X}, \vec{x})$ depends on the choice of external and internal coordinates. Considering a system that is symmetric with respect to the external coordinates, the integration over the external degrees of freedom gives a factor of V for the integration volume and $8\pi^2$ for the rotations:

$$Q = 8\pi^2 V \xi \int d\vec{x} J(\vec{x}) e^{-\beta \bar{E}(\vec{x})} \quad (2.51)$$

Here $\bar{E}(\vec{x})$ is regarded as the potential of mean force with the degrees of freedom of the solvent averaged out. Considering now a system with a receptor and a ligand. The partition function of both the receptor and the ligand have the form

$$Q = 8\pi^2 V \xi \int d\vec{x}_i J_i(\vec{x}_i) e^{-\beta \bar{E}_i(\vec{x}_i)}, \quad (2.52)$$

with $i = R, L$ for the receptor and the ligand. The energy function of the ligand and the receptor can be written as the sum of the energies of the ligand E_L and the receptor E_R and an interaction term w that depends on the distance \vec{r} and the relative orientation $\vec{\omega}$ of the receptor and the ligand.

$$\bar{E}_{RL}(\vec{x}_{RL}) = \bar{E}_R(\vec{x}_R) + \bar{E}_L(\vec{x}_L) + w(\vec{r}, \vec{\omega}) \quad (2.53)$$

Then by inserting the partition function into eq. 2.41 and using the notation $\vec{x}_{RL} = (\vec{x}_R, \vec{x}_L, \vec{r}, \vec{\omega})$ the association constant can be calculated by

$$K_a = \frac{(8\pi^2)^{-1} \int_b d\vec{x}_{RL} J_{RL}(\vec{x}_{RL}) e^{-\beta \bar{E}_{RL}(\vec{x}_{RL})}}{\int d\vec{x}_R J_R(\vec{x}_R) e^{-\beta \bar{E}_R(\vec{x}_R)} \int d\vec{x}_L J_L(\vec{x}_L) e^{-\beta \bar{E}_L(\vec{x}_L)}}, \quad (2.54)$$

where the integral over the partition function of the complex Q_{RL} has to be restricted to the region in the configurational space where the complex is formed denoted by b [140]. It is possible to define a potential of mean force $W(\vec{r}, \vec{\omega})$ by

$$e^{-\beta W(\vec{r}, \vec{\omega})} := \frac{\int d\vec{x}_R d\vec{x}_L J_R(\vec{x}_R) J_L(\vec{x}_L) e^{-\beta \bar{E}_{RL}(\vec{x}_R, \vec{x}_L, \vec{r}, \vec{\omega})}}{\int d\vec{x}_R J_R(\vec{x}_R) e^{-\beta \bar{E}_R(\vec{x}_R)} \int d\vec{x}_L J_L(\vec{x}_L) e^{-\beta \bar{E}_L(\vec{x}_L)}} \quad (2.55)$$

Then the binding constant K_a can be obtained as:

$$K_a = (8\pi^2)^{-1} \int d\vec{r} d\vec{\omega} J_r(\vec{r}) J_\omega(\vec{\omega}) e^{-\beta W(\vec{r}, \vec{\omega})} \quad (2.56)$$

By integration over the rotational degrees of freedom in $\vec{\omega}$ a potential of mean force $\bar{W}(\vec{r})$ can be defined that only depends on the distance vector \vec{r} :

$$e^{-\beta \bar{W}(\vec{r})} := (8\pi^2)^{-1} \int d\vec{\omega} J_\omega(\vec{\omega}) e^{-\beta W(\vec{r}, \vec{\omega})} \quad (2.57)$$

With $\bar{W}(\vec{r})$, the binding constant K_a can be obtained as:

$$K_a = \int_b d\vec{r} e^{-\beta \bar{W}(\vec{r})} \quad (2.58)$$

The potential of mean force $\bar{W}(\vec{r})$ contains enthalpic and entropic contributions of whole configuration space of the ligand and the receptor. Substituting K_a in eq. 2.43 results in

$$e^{-\beta\Delta G_b^0} = \int_b d\vec{r} e^{-\beta\bar{W}(\vec{r})} \quad (2.59)$$

Therefore, when $\bar{W}(\vec{r})$ can be approximated, it allows the calculation of the standard binding free energy ΔG_b^0 without calculating the partition functions Q_i . This fact is exploited in simulation techniques which integrate the approximated potential of mean force (e.g. umbrella sampling [120]).

2.5 Molecular Dynamics Simulations

Biological systems are often complex many-particle systems for which, contrary to crystalline or solid state systems, no straightforward reduction to a few degrees of freedom is possible. The temperature of interest for these systems typically lies close to 300 K. Therefore, entropic effects significantly contribute to the thermodynamics of these systems, requiring the explicit consideration of many degrees of freedom in order to adequately describe their state [125]. In the last fifty years, structural biology has provided atomic-resolution models of many molecules that are essential to life, including proteins, nucleic and ribonucleic acids. The static molecular structures determined by X-ray crystallography and other techniques are tremendously useful, but in reality, the molecules they represent are highly flexible and their dynamics are often critical to their function. Proteins, for example, undergo a number of conformational changes during their lifetime. Important examples are the folding of their secondary, tertiary and quaternary structures or functional conformational changes that enable proteins to e.g. catalyse reactions, transport molecules, transduce signals, etc. The dynamics of such processes is sometimes crucial for the understanding of a biomolecule. Experimental techniques that provide information about the dynamics of biomolecules are generally limited in their spatial and temporal resolution and often provide information about the time and/or ensemble averages of molecules instead of individual molecules. An alternative to experimental observation is modeling that uses the knowledge about the atomic structure and the underlying physical laws to compute the dynamics of the molecule [32]. Usually, the electronic degrees of freedom are the highest level of theory that is considered in the modeling of biological systems. The dynamics of such a system is described by the time-dependent Schrödinger equation:

$$\mathcal{H}\psi(\vec{R}, \vec{r}, t) = i\hbar \frac{\partial \psi(\vec{R}, \vec{r}, t)}{\partial t} \quad (2.60)$$

where \hbar is the reduced Planck constant, i the imaginary unit, \mathcal{H} denotes the Hamilton operator of the system and $\psi(\vec{R}, \vec{r}, t)$ the wave function of the system as a function of the coordinates of the nuclei \vec{R} , the electrons \vec{r} at the time t . However, the system sizes of biological macromolecules and the timescales of interest render the explicit solution of eq. 2.60 computationally unfeasible.

The standard method for the calculation of macromolecular dynamics is known as all-atom molecular dynamics (MD) simulations, where a representation of the biological system evolves in time according to the laws of classical mechanics. The representation is based on position, velocities of and forces between particles which represent one or more atoms. Sometimes more than one particle describe a single atom. For example, the oxygen in the TIP4P water model [67] is described by 2 virtual particles needed for the accurate description of the electrostatics. In the following, I describe the major approximations applied in MD simulations that allow simulation of the molecular dynamics of biological systems at atomic resolution.

2.5.1 Approximations

The first approximation from an *ab initio* treatment is the separation of the electronic and the nucleic degrees of freedom according to the *Born-Oppenheimer approximation*. The fraction of the mass of an electron m and a nucleus M is typically

$$\frac{m}{M} \approx 10^{-3} - 10^{-5} \quad (2.61)$$

Therefore, the electrons move much faster than the nuclei and it is possible to approximate the potential of the nuclei as quasistatic. Assuming that the wave function of the electrons changes adiabatically with the potential generated by the nuclei the wave function of the system can be expressed as a product of an electronic ψ_{el} and a nucleic wave function ψ_{nuc} :

$$\psi(\vec{R}, \vec{r}, t) = \psi_{\text{nuc}}(\vec{R}, t) \cdot \psi_{\text{el}, \vec{R}}(\vec{r}) \quad (2.62)$$

where the electronic wave function depends on the nucleic coordinates only parameterically. This approximation is known as the *Born-Oppenheimer approximation*. When \mathcal{H} is the Hamilton operator of the complete system, and T_{nuc} the kinetic energy contribution of the nuclei, the electronic Hamilton-Operator can be defined as $\mathcal{H}_{\text{el}} := \mathcal{H} - T_{\text{nuc}}$. For a given nucleic configuration \vec{R} , the application of \mathcal{H}_{el} on the electronic wave function gives the eigenvalue $E_{\text{el}}(\vec{R})$ or

$$\mathcal{H}_{\text{el}}\psi_{\text{el}}(\vec{r}, \vec{R}) = E_{\text{el}}(\vec{R}) \cdot \psi_{\text{el}, \vec{R}}(\vec{r}) \quad (2.63)$$

When the system of interest is treated in the electronic ground state the smallest eigenvalue $E_{\text{el}}^0(\vec{R})$ is used and the time dependent Schrödinger equation (eq. 2.60) can be rewritten as

$$[T_{\text{nuc}} - E_{\text{el}}^0(\vec{R})] \psi_{\text{nuc}}(\vec{R}, t) = i\hbar \frac{\partial \psi_{\text{nuc}}(\vec{R}, t)}{\partial t} \quad (2.64)$$

Since the potential energy is now independent of time, the wave function of the nuclei can be separated in a time dependent part

$$\phi(t) = e^{-iEt/\hbar} \quad (2.65)$$

and a time independent part $\psi(\vec{R})$. Replacing $\psi_{\text{nuc}}(\vec{R}, t)$ in equation 2.64 by the product $\psi(\vec{R}) \cdot \phi(t)$, applying the time derivative and multiplying by ψ^* results in an equation for the dynamics of the nuclei:

$$[T_{\text{nuc}} - E_{\text{el}}^0(\vec{R})] \psi_{\text{nuc}}(\vec{R}) = E \psi_{\text{nuc}}(\vec{R}) \quad (2.66)$$

Here E is the total energy of the system. The potential $E_{\text{el}}^0(\vec{R})$ also contains the coulomb interaction energy between the nuclei.

Parameterization of the Electronic Potential

The calculation of the effective electronic potential $E_{\text{el}}^0(\vec{R})$ requires the solving of the time-dependent electronic part (Eq. 2.63). At present, it is not feasible to solve this equation for more than a few atoms. Thus, the potential is approximated by a classical molecular mechanics (MM) energy function U :

$$E \approx U = \sum_{\text{bond}, i} U_{\text{b}}^i + \sum_{\text{angles}, j} U_{\text{a}}^j + \sum_{\text{dihedrals}, k} U_{\text{d}}^k + \sum_{\text{imp}, l} U_{\text{i}}^l + \sum_{\text{pairs}, m, n} (U_{\text{vdw}}^{m, n} + U_{\text{coul}}^{m, n}) \quad (2.67)$$

which is the second major approximation applied in MD simulations. As shown by eq. 2.67, a typical MM energy function for bio-molecular applications consists of individual *bonded* potential terms for bonds, bond-angles, dihedrals and improper dihedrals describing the forces of the covalent bonds and *non-bonded* potentials that describe London-dispersion, Pauli-repulsion and Coulomb interactions. These functions are kept relatively simple, because of the need to evaluate these functions a large number of times during the simulation. The form of U and the parameters used in U are referred to as the *force field*. Unfortunately, the term force field is not homogenously used and can mean either the system of potential energy functions or the set of force constants, i.e. the partial derivatives of U according to the atomic coordinates. In general, the force field

contains many adjustable parameters which are often obtained by fitting to data from experiments or from *ab initio* calculations. The majority of contemporary force fields keep the partial charge of the particles fixed, and therefore do not explicitly account for electronic polarizability.

Newtonian Dynamics

The third major approximation is the assumption that the nuclei can be modelled by a set of point masses which follow the *newtonian equations of motion*:

$$m_i \frac{d^2 \vec{R}_i(t)}{dt^2} = \nabla_i U(\vec{R}_1 \dots \vec{R}_n) = \vec{F}_i \quad (2.68)$$

For numeric integration of these equations the leap-frog integrator,

$$\vec{v}(t + \frac{1}{2}\Delta t) = \vec{v}(t - \frac{1}{2}\Delta t) + \frac{\Delta t}{m} \vec{F}(t) \quad (2.69)$$

$$\vec{r}(t + \Delta t) = \vec{r}(t) + \Delta t \cdot \vec{v}(t + \frac{1}{2}\Delta t) \quad (2.70)$$

implemented in Gromacs [53, 124] was used for all molecular dynamics simulations. Here $\vec{r}(t)$ are the atomic coordinates and $\vec{v}(t)$ the velocities. To obtain a ensemble of constant pressure the Parrinell-Rahman [98] barostat (Section 2.5.2) was applied, and to obtain an ensemble of constant temperature the v-rescale thermostat [18] (Section 2.5.2) was applied. In this framework, these equations of motion are modified to couple the box size to the microscopic motion of the particles.

2.5.2 General Simulation Conditions

All MD simulations presented herein were performed utilizing the Gromacs simulation package [53, 124]. Canonical isothermal-isobaric (NPT) ensembles were generated by coupling the simulation box to a target temperature T_t and target pressure p_t , keeping the number of particles, the pressure and the temperature constant. The coupling of the system is described in the following paragraphs, as well as the force field that was used (Amber ff99SB-ILDN).

Constant Temperature

Controlling the temperature is an important issue in MD simulations. For the simulations in chapter 4 the v-rescale thermostat [18] was used that is an extension to the Berendsen thermostat [8]. The v-rescale thermostat uses a stochastic

term to generate a proper canonical ensemble. It was derived as a modification of the standard velocity rescaling that, at a predefined frequency, scales all velocities by factor α

$$\alpha = \sqrt{\frac{E_k^t}{E_k}}, \quad (2.71)$$

with the kinetic energy of the system E_k and the target kinetic energy E_k^t that corresponds to a target temperature T_t . In order to ensure a canonical ensemble, E_k^t is drawn from a canonical equilibrium distribution. This demonstrates the principle of the v-rescale thermostat; the actual implementation is a bit different:

$$dE_k = (E_k^t - E_k) \frac{dt}{\tau} + 2 \sqrt{\frac{E_k^t E_k}{N_f}} \frac{dW}{\sqrt{\tau}} \quad (2.72)$$

Here the applied change in the kinetic free energy dE_k depends on the current value of E_k . N_f is the number of degrees of freedom. This keeps the changes of the velocities small and results in a relatively undisturbed system. The parameter τ controls the strength of the coupling. The term dW generates noise that is needed to generate a valid thermodynamic ensemble.

Constant Pressure

The Parrinello-Rahman [98] barostat implemented in Gromcas [53, 124] was used to keep the pressure constant in the simulations reported in this work. The general idea behind the barostat is to couple the microscopic motions of the particles to the pressure. Here I will, briefly sketch the underlying mathematics. Using the Clausius virial theorem the pressure tensor \mathbf{P} can be obtained as:

$$\mathbf{P} = \frac{2}{V} (E_k - \Gamma), \quad (2.73)$$

with the inner virial tensor

$$\Gamma = -\frac{1}{2} \sum_{i < j} \vec{r}_{ij} \cdot \vec{F}_{ij} \quad (2.74)$$

The isotropic pressure p results from the trace of \mathbf{P}

$$p = Tr(\mathbf{P})/3 \quad (2.75)$$

Correcting the pressure (also anisotropically) in a simulation can be achieved through a change in the inner virial Γ by scaling the inter particle distances. The pressure is then kept constant by coupling the box size to the microscopic motion of the atoms.

The Amber99sb Force-Field

All simulations that are were reported in this work were performed with the Amber ff99SB [57] and the Amber ff99SB-ILDN force fields. Some characteristic features of these force fields are the use of fixed partial charges on atom centers, the explicit use of all hydrogen atoms and no specific functional forms for hydrogen bonding. The protein ϕ/ψ dihedrals have specific rotational parameters that affect the relative energies of alternate backbone conformations. These parameters were fit to relative energies obtained from *ab initio* calculations of alternate rotamers of small molecules. The partial atomic charges were derived by fitting the electrostatic potential calculated with a Hartree-Fock method using the 6-31G* basis set in the gas phase. This approach was chosen because it overpolarizes bond dipoles, such that the resulting charge distribution better approximates those occurring in aqueous phase. The calculated potential is used to fit point charges at the center of each atom in a way that electrostatic potential of the point charges mimics the calculated potential at the surface of the molecules. The original force field Amber ff99 [132] (which was a modification of the Amber ff94) was known to overestimate the α -helical phase. This overestimation was resolved by a set of new ϕ/ψ parameters in the ff99SB extension[80]. The Amber ff99SB-ILDN force field has slightly modified torsion parameters for Ile (I), Leu (L), Asp (D) and Asn (N) side chains. The improvement was validated by comparison data from a microsecond-timescale MD simulation to nuclear magnetic resonance (NMR) data that probed these torsions directly. The newest force field of the Amber family is the ff99SB*-ILDN [11], that applies further modifications to the ϕ/ψ -angles. This force field was not implemented in Gromacs by the time of my studies and therefore was not used.

2.5.3 Limitations

Limitations of the MD simulations arise from the approximations described above and from computational deficiencies which limit the scope of feasible timescales and system sizes. At first, I will discuss the basic limitations, which arise from the approximations mentioned above. Afterwards, I will focus on practical limitations and cover the state of the art of contemporary MD simulations.

Only What the Code Allows

The most basic level limitations lie in the simulation algorithm itself. A simulation can only generate results which lie in the space that can be covered by the algorithm that is used. For example, a MD simulation that uses an integration

time step of 1 fs can not generate vibrations with a frequency above 5 PHz. The next basic limitations are caused by the approximations in the construction of the model. Effects which are completely neglected in the model building will not appear in the simulation. For example, a simulated $^{235}_{94}\text{P}$ atom will not undergo a nuclear reaction, since effects of the weak interaction are neglected in the standard framework of MD simulations. A more critical restriction is the conservation of covalent bonds that cannot break or be formed in a standard MD simulation. Chemical reactions are currently not part of the parameterization process and cannot be simulated with pure MD simulations.

However, there are phenomena which are not explicitly covered in the model building step, whose effects are fed implicitly into the simulation. This holds, for example, for hydrogen bonds or mutual polarization, effects which are based on changes of electronic wave functions. An example is the overestimation of the molecular partial charges as mentioned in the next paragraph. To some extent, the effects of these phenomena are included in an MD simulation, since, in general, force field parameters are fitted against *ab initio* and/or experimental findings. The fitting process therefore partially compensates inadequacies of the basic model.

Force-Fields are Fitted

The parameterization process determines which properties can be reliably determined from an MD simulation. A particular force field that is trained to reproduce correct thermodynamics can be inaccurate in the description of transition kinetics. One force field may be parameterized in order to obtain correct diffusion constants and others in order to predict relative conformational energies. Therefore, the choice of which force field to use depends on the properties of interest and the question that is addressed.

Also the timescales that were used in the parameterization process of the force field play a role. On larger timescales inadequacies of the force field may become critical, for example the overestimation of α -helical arrangements. In a recent study, a number of force fields was evaluated at the microsecond timescale [79]. The results indicated that the examined force fields have consistently improved over the last ten years. The force fields with the best agreement with the experimental data were Amber ff99SB*-ILDN [11] and CHARMM22*. The most recent versions provide an accurate description of many dynamic and structural properties of the tested proteins. However, none of the force fields was able to accurately capture the temperature dependency of the secondary structures. The same most likely holds for larger system sizes.

Maximal Scales

One of the largest biological systems described by MD simulations was the southern bean mosaic virus with 4.5 million atoms [143, 144]. In practice MD simulations were until recently generally limited to nanosecond timescales. However, recent advances in hardware and algorithms have increased the timescales accessible to MD simulations to the millisecond-scale, allowing MD to capture critical biochemical processes that take place on these timescales such as protein folding [79], protein ligand binding [17, 32, 114], and major conformational changes which are essential to protein function [63].

2.5.4 Calculation of the Binding Free Energy

As mentioned in section 2.4.3, for the condition of constant pressure, temperature and number of particles, the Gibbs free energy $G(T, p, N)$ is a thermodynamic potential that describes the stability of a state corresponding to the natural variables of G : the temperature T , the pressure p and the number of particles N . The difference of the free energy $\Delta G_{A,B}$ between two states A and B describes the relative stability of the states. The gradient of the free energy $\nabla G_{A,B}$ with respect to the natural variables describes the driving forces of thermodynamic system. In equilibrium, G is minimized and $\nabla G_{A,B}$ becomes zero. Considering the complexation of a ligand and a receptor, the knowledge about difference of the Gibbs free energy $\Delta G_{A,B}$ between the dissociated state A and the complexed state B allows the deduction of essential properties relevant for ligand receptor binding as e.g. the equilibrium constant K_a (Section 2.1.1 and 2.4.3), the IC_{50} as well as enthalpic and entropic contributions of binding (Section 2.4.4).

According to equations 2.41 and 2.29 it could be argued that the exact calculation of ΔG is in principle possible. However, the large number of conformational states a macromolecule can sample renders an exact calculation of Eq. 2.29 practically unfeasible. Equation 2.9 gives an idea about this number. Therefore, the calculation of the partition function Q is only possible for very simple systems, but in general impracticable for highly complex systems such as solvated macromolecules. However, not all parts in phase space have the same impact on the absolute values of Q . Therefore it is possible to take structural ensembles generated with simulation techniques for the effectual approximation of Q , here referred to as the effective partition function \bar{Q} . This fact is used, for example, in the methods that base on MD simulations.

ΔG or differences of ΔG s, denoted as $\Delta\Delta G$, can be approximated from MD simulations. The *absolute binding free energy* is the free energy difference between

the dissociated state and the complexed state of a ligand and a receptor. The *relative binding free energy* $\Delta\Delta G$ is the difference between the absolute binding free energies of two complexes. This can be related to the thermodynamic cycle in Fig. 2.5. ΔG_2 and ΔG_4 are the absolute binding free energies of the ligands L_2 and L_1 which can be measured experimentally. The transitions along the horizontal arrows correspond to modification of the ligands from L_1 to L_2 in solvent (ΔG_1) and in the complex (ΔG_3), which are impossible to measure experimentally.

Different molecular dynamics (MD) or monte carlo based simulation techniques can be used to calculate the transitions corresponding either to the horizontal or the vertical arrows. The vertical arrows correspond to different regions in phase space and are described by the same Hamilton operator. When the states A and B can be separated into two non-overlapping regions in phase space it is, in principle, possible to use standard MD simulations to obtain the relative weights (Q_A and Q_B as described in section 2.4.2). The free energy between these states is related to these weights by

$$\Delta G_4 \approx -k_B T \ln \left(\frac{\bar{Q}_A}{\bar{Q}_B} \right) \quad (2.76)$$

If an MD simulation samples transitions between the states A and B the weights can be estimated by the probabilities to find the system in either state A or B . Unfortunately, this very direct method has practical limitations. The transition timescale has to be small enough to be crossed frequently by an MD simulation in order to reach convergence. This requires both the free energy difference and the barrier between both states to be small. Ligand-receptor associations can easily last several microseconds and the dissociation can be even slower by orders of magnitude [17, 32, 114]. This renders the calculation of ΔG_2 and ΔG_4 by free MD simulations practically impossible for high affinity ligands.

Umbrella Sampling

This limitation can be overcome by applying enhanced sampling methods as, for example, umbrella sampling (US). The basic idea behind the US is that the relative probability to sample a point \vec{r} in configuration space is not given by

$$p(\vec{r}) = e^{-\beta\bar{W}(\vec{r})} \quad (2.77)$$

but

$$p_r(\vec{r}) = w(\vec{r})e^{-\beta\bar{W}(\vec{r})}, \quad (2.78)$$

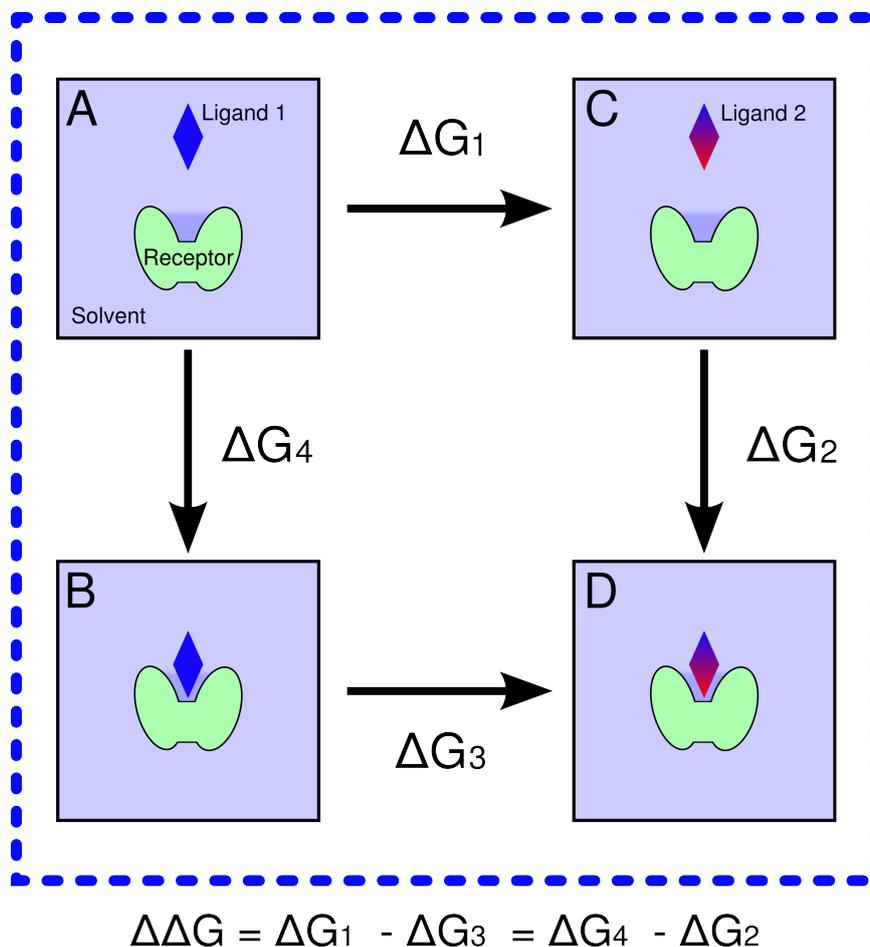


Figure 2.5: Thermodynamic cycle for the calculation of the relative binding free energy of receptor-ligand complexes with ligands L_1 and L_2 $\Delta\Delta G$ (blue cycle): **A:** Receptor and dissociated ligand L_1 in solvent. **B:** Solvated complex with ligand L_1 . **C:** Receptor and dissociated ligand L_2 in solvent. **D:** Solvated complex with ligand L_2 . By routing the system to different states along the arrows the corresponding free energy differences are attained. The binding free energies of complexation of a ligand with a receptor ΔG_2 and ΔG_4 (vertical arrows) can be measured experimentally or calculated with the umbrella sampling method. The horizontal arrows represent the modification of a ligand L_1 to another structure L_2 . The corresponding free energy differences ΔG_1 and ΔG_3 can be calculated using thermodynamic integration or free energy perturbation. If L_2 consists of dummy atoms without physical interactions with the environment ΔG_2 becomes zero.

where $w(\vec{r})$ is a weighting function and $\bar{W}(\vec{r})$ the potential of mean force (PMF). This can be obtained by adding an umbrella potential [120] $U(\vec{r}, \vec{r}_0)$ to the Hamilton operator that restrains the system to the point \vec{r}_0 . For example, a harmonic potential:

$$U(\vec{r}, \vec{r}_0) = \frac{1}{2} \kappa (\vec{r} - \vec{r}_0)^2, \quad (2.79)$$

where $\kappa > 0$ is the force constant. The unrestrained probability $p(\vec{r})$ for the point \vec{r} can be expressed in terms of the restrained probability $p_r(\vec{r})$ and an ensemble average $\langle \dots \rangle_r$ generated with the Hamilton operator $\mathcal{H}_r(\vec{r})$:

$$p(\vec{r}) = \frac{p_r(\vec{r}) e^{\beta U(\vec{r}, \vec{r}_0)}}{\langle e^{\beta U(\vec{r}, \vec{r}_0)} \rangle_r}, \quad (2.80)$$

which is the desired probability. By routing the system with the umbrella potential from state A to state B it is therefore possible to calculate ΔG_2 (ΔG_4 analogously). In practice, $p(\vec{r})$ is calculated for many points \vec{r}_i along a path that combines the phase space regions of states A and B . Then the underlying PMF is calculated using e.g. the weighted histogram analysis method (WHAM) [77].

Alchemical Methods

Alternatively, the free energy differences corresponding to the horizontal arrows in Fig. 2.5 can be calculated by using the fundamental equation of free energy perturbation known as the Zwanzig [148] equation:

$$\Delta G_{A,B} = G_B - G_A = -k_B T \ln \left(\langle e^{-\beta(\mathcal{H}_B - \mathcal{H}_A)} \rangle_A \right) \quad (2.81)$$

Here \mathcal{H}_i represents the Hamilton operator of the i th state and $\langle \rangle_i$ an ensemble average of the system described with the Hamilton operator \mathcal{H}_i . Here, the states A and B are not different regions in configuration space, but systems with different energy functions \mathcal{H}_A and \mathcal{H}_B . The idea behind Eq. 2.81 is that, when state A and state B significantly overlap in phase space the free energy of state B can be obtained from the ensemble generated with Hamiltonian \mathcal{H}_A from state A as:

$$G_B \approx \int_A d\vec{x} e^{-\beta \mathcal{H}_B(\vec{x})}, \quad (2.82)$$

where \vec{x} is a point in phase space that corresponds to state A .

Perturbation approaches generalize Eq. 2.81 to

$$\Delta G_{A,C} = \sum_{i=0}^{n-1} \Delta G(\lambda_i, \lambda_{i+1}) = \sum_{i=0}^{n-1} -k_B T \ln \left(\langle e^{-\beta(\mathcal{H}_{\lambda_{i+1}} - \mathcal{H}_{\lambda_i})} \rangle_{\lambda_i} \right), \quad (2.83)$$

where states λ_0 and λ_n correspond to state A and C in Fig. 2.5. The difference between states λ_i and λ_{i+1} is obtained by evaluating the ensemble generated with

the Hamilton operator \mathcal{H}_i corresponding to λ_i . In analogy to Eq. 2.83, $\Delta G_{C,A}$ is defined by

$$\Delta G_{C,A} = \sum_{i=0}^{n-1} \Delta G(\lambda_{n-i}, \lambda_{n-i-1}) \quad (2.84)$$

$\Delta G_{A,C}$ is referred to as the *forward perturbation* and $\Delta G_{C,A}$ to as the *backward perturbation*. The convergence of $\Delta G_{A,C}$ and $-\Delta G_{C,A}$ is necessary but not sufficient for an accurate calculation of ΔG_1 , since the transitions can be performed too fast for the system to equilibrate. Based on the works of Jarzynski [61] and Crooks [23], in recent years, the non-equilibrium techniques in combination with the Bennett's acceptance ratio method (BAR) have become the state of the art, as they have shown to be computationally more efficient than FEP [16].

A frequently reported technique for the calculation of free energy differences is thermodynamic integration (TI) [71]. Here, the Hamilton operator describing the system is routed from state A to state B , by a parameter λ . The Hamilton operator

$$H(\lambda) = \lambda \mathcal{H}_B + (1 - \lambda) \mathcal{H}_A \quad (2.85)$$

and its derivative with respect to λ

$$\mathcal{H}'(\lambda) = \frac{\partial \mathcal{H}(\lambda)}{\partial \lambda} \quad (2.86)$$

are introduced. The parameter λ is treated as a reaction coordinate that routes the system from state A to state B by going from 0 to 1. The free energy difference $\Delta G_{A,B}$ is then calculated by performing the integration:

$$\Delta G = \int_{\lambda=0}^{\lambda=1} d\lambda \langle \mathcal{H}'(\lambda) \rangle_{\lambda} \approx \sum_{i=1}^n \Delta \lambda \langle \mathcal{H}' \rangle_n \quad (2.87)$$

A similar approach, discrete TI, uses discrete steps of λ and to evaluate $\Delta G_{A,B}$ at intermediate states [74] as indicated by the inequality in Eq. 2.87. The intermediate states λ_i with $i = 1, \dots, n-1$ often correspond to unphysical states. Therefore, these methods are called *alchemical*.

Summary

Both types of free energy calculations (enhanced conformational sampling and alchemical calculations) require converged ensemble averages for all intermediated states r_i or λ_i . Therefore, considerable computational effort is required to perform free energy calculations for the estimation of ligand-receptor binding free energies [93]. Currently, this effort limits the application of such calculations

practically to a few (1–100) compounds. Especially, the calculation of absolute binding free energies is computationally costly, since a larger perturbation, namely the complete (dis)appearance of a compound needs to be simulated.

3

Optimization of Molecular Docking

The current chapter provides a case study for the design, optimization and application of a combined virtual screening (VS) approach, that was experimentally validated using an automated patch clamp technique. Initially, four different molecular docking programs (Autodock-Vina, eHiTS, FlexX and Glide) were benchmarked against a training set of 2576 compounds that were provided by Xention (<http://www.xention.com>). Here, computations with Glide were performed by Dr. Wiktor Jurkowski at the Department of Biochemistry and Biophysics of the Stockholm University (Sweden); and computations with eHiTS were performed by Dr. Katie J. Simmons at the Chemistry Department of the University of Leeds (United Kingdom). The reported automated patch clamp experiments were performed by Dr. Jean-Francois Rolland at Xention (Cambridge, United-Kingdom).

3.1 Introduction

Potassium channels are a diverse protein class with at least 78 different members. They virtually appear in all living organisms and are necessary for many physiological processes including cell excitability and secretion mechanisms. The voltage-dependent potassium channels are the largest sub-group of potassium channels with 12 families of phylogenetically related proteins (K_V1-12). They

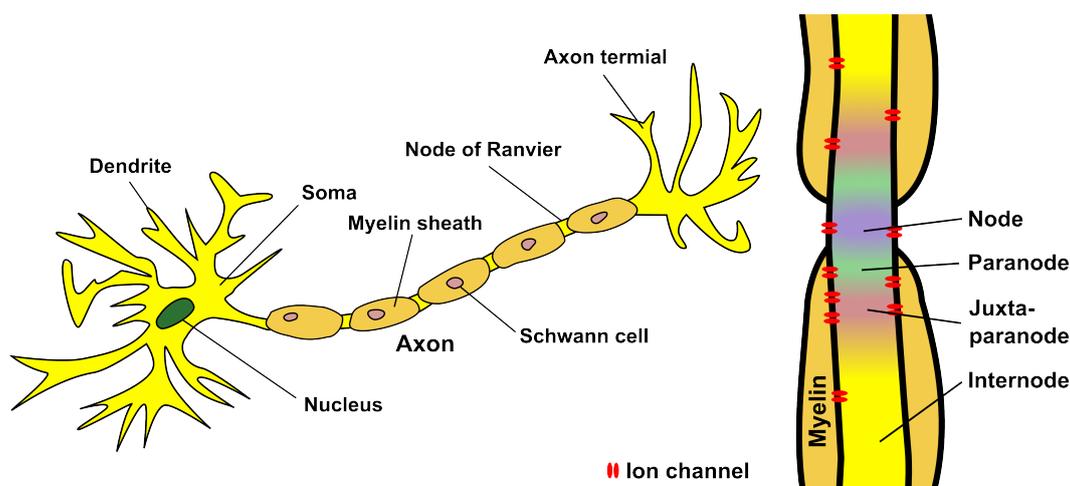


Figure 3.1: Scheme of a typical neuron and illustration of the Nodes of Ranvier and the internodal, juxtapanodal and paranodal regions.

appear in virtually all living organisms and play a crucial role in many life-sustaining functions including neural transduction and the regulation of the heart beat frequency [69]. Accordingly, their dysfunction accounts for a wide range of human pathologies as episodic ataxia, atrial fibrillation or the long QT syndrome. Some of the K_V1 members form heteromultimers, with the $K_V1.1$ and $K_V1.2$ combination being one of the most abundant in both the central and peripheral nervous system (CNS and PNS, respectively) [21, 102, 131]. More importantly, this type of heteromers has been found to specifically co-localize at myelin-protected juxtapanodal regions of the nodes of Ranvier of nerve axons (Figure 3.1) [3, 29, 103, 130], where they control axon excitability and ensure saltatory conduction [109]. In demyelinating diseases such as multiple sclerosis (MS), these channels are exposed and nerve conduction is impaired [69]. Recently, 4-aminopyridine (INN:Fampridine), a non-selective potassium channel inhibitor, has been approved as the first medication to improve the ability to walk in people suffering from multiple sclerosis (MS) [46, 47, 48], probably by blocking the exposed $K_V1.1-1.2$ heteromeric potassium channels. (A list of known $K_V1.1$ or $K_V1.2$ inhibitors and references is provided in supporting information Tab. 7.2). Unfortunately, its low potency and poor channel specificity raise issues, particularly in regard to cardiac safety. Therefore, the search for more selective blockers, and the development of proper strategies for the study of drug-channel interactions, are highly desirable from a clinical perspective.

$K_V1.1$ and $K_V1.2$ are tetrameric potassium channels with a single pore in the center that is selectively permeable for water and potassium (Fig. 3.2). Each subunit contains a transmembrane domain that is composed of six membrane

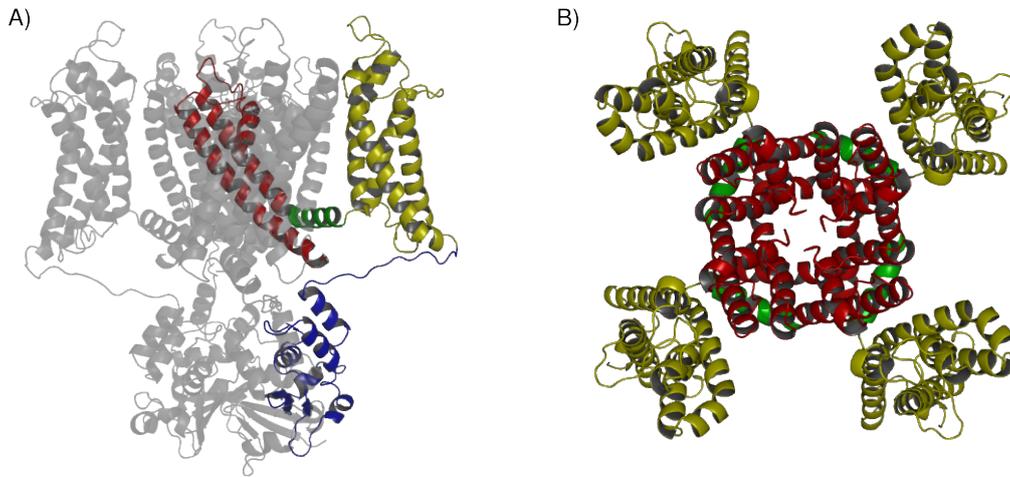


Figure 3.2: Ribbon representation of the voltage gated potassium channel $K_V1.2$. Lateral view (A) and view from the exoplasmic side (B). The structure is colored according to individual functional components: pore domain (red), linker (green), voltage sensing domains (yellow) and the N-terminal domain (blue) which is not depicted in the top view.

spanning α -helices. Four of these helices form a separate voltage sensing domain (the β -subunit) whereas the remaining two helices are part of the central pore (α -subunit), whose sequence is almost conserved in this family of potassium channels. Both sub-units are linked by a small α -helix, termed the linker. Per sub-unit, this domain consists of two TM- α -helices connected by a loop. These loops connecting the transmembrane helices of the α -sub-units contain the key element for the ion selectivity decoded in a highly conserved motif of eight amino acids (TXXTXGXG). This sequence serves as signature for the identification of possible potassium channel coding genes [51, 108]. Five of these amino acids (TVGYG in $K_V1.2$) are situated on an elongated region of the pore loop near the extracellular vestibule of the channel 3.3. This arrangement is also conserved in other potassium channels [31, 64, 78, 85, 141] and is termed the *selectivity filter* (Figure 3.3). The selectivity filter can adopt a conformation where it is able to conduct potassium ions. In this conformation the amino acids which form the selectivity filter are arranged in a way that the backbone peptide carbonyl groups point towards the center of the tetramers. Here the carbonyl oxygen atoms form a ring mimicking the coordination of the hydration shell oxygen atoms of a potassium ion in water. Such an arrangement has been shown to allow the dehydration of small cations, which are strongly bound to water molecules in bulk solution, with a reduced free energy barrier [141].

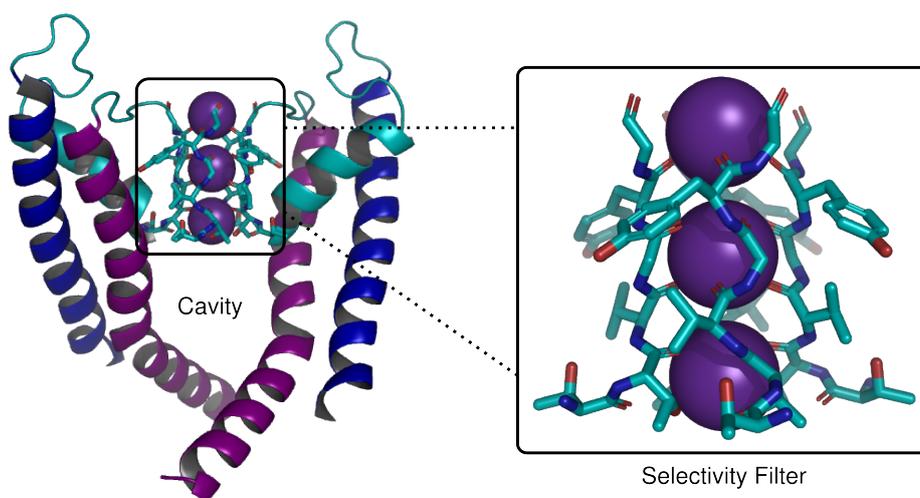


Figure 3.3: Structure of the α -subunits (left). For clarity, only two chains are shown. The residues of the selectivity filter are shown in stick representation (right). The purple spheres illustrate K^+ ions in the filter.

Importantly, $K_V1.2$ is among the few channels for which structural insight has been gained through X-ray crystallography [84, 85]. This enables to perform structure based virtual screening. However, most docking schemes have been developed to work on well-defined binding pockets such as enzymatic active sites, where docking poses can be compared with specific pharmacophores. How well these programs enrich active compounds specifically on the putative binding sites such as the inner cavity of a potassium channel remains unclear, since no VS benchmark regarding potassium channels has been reported. In studies of KcsA and $K_V1.5$, the inner cavities were used successfully as target sites in structure based virtual screening (SBVS) [82, 83]. Liu *et al.* [82] screened 200,000 molecules from the Accelrys Available Chemicals Directory (<http://accelrys.com>) against the extracellular pore entrance of the pH gated potassium channel KcsA. The 300 top scored molecules were then optimized in a molecular mechanics force field and the interaction energies were calculated. According to the docking scores, the calculated interaction energies and calculated solubility coefficients (LogP values), 20 compounds were selected and assayed *in vitro*. Six compounds were found to suppress the K^+ conductance by more than 10% at a compound concentration of $100 \mu\text{M}$ when applied from the extracellular site. Yang *et al.* [138, 139] screened the Maybridge (<http://www.maybridge.com>) library against a homology model of $K_V1.5$. The top 1000 compounds were evaluated in a consensus molecular docking approach. Then, 18 compounds were manually selected and assayed. Five of the 18 compounds blocked $K_V1.5$ mediated ion

current by more than 50 % at a compound concentration of 10 μ M. Furthermore, modeling and mutagenesis studies confirm that the inner cavity is a binding site for small ionic molecules [2, 27, 86]. Eldstrom and Fedida [34] used Autodock4 to model the high affinity binder Vernakalant into the cavity of a homology model of $K_V1.5$. The manual selection that was used in the virtual screening studies renders the estimation of the hit rate of the VS approaches difficult.

In this collaborative study, four popular molecular docking approaches (eHiTS, FlexX, Glide, and Autodock-Vina) were benchmarked for their ability to distinguish between compounds known to be active or inactive against the potassium channel concatamer $K_V1.1-(1.2)_3$, consisting of one $K_V1.1$ subunit and three $K_V1.2$ subunits. An effective virtual screening technique was established. For that purpose, a commonly used strategy for the improvement of molecular docking predictions was applied, namely consensus scoring (CS) (Section 2.3.4). The aim of the current study was to benchmark, optimize, and employ molecular docking based techniques for a specific target ($K_V1.1-(1.2)_3$) in order to find novel and potent inhibitors, thereby decreasing both experimental effort and costs. Moreover, to eliminate any cardiac liability, all identified hits were assayed against cardiac channels known to be involved in the cardiac action potential.

3.2 Results

The quality of four popular molecular docking approaches (eHiTS, FlexX, Glide, and Autodock-Vina) was assessed by screening a library of compounds known to be active or inactive on $K_V1.1-(1.2)_3$. Owing to the total lack of pharmacological modulators of the $K_V1.1-(1.2)_3$, we mined an in-house compound library from the pharmaceutical company Xention (<http://www.xention.com>). 2675 active and inactive compounds were selected to set up a test-library. The fraction of active compounds in this test-library was 32%. After the initial screen, the most predictive terms of the individual docking procedure were combined to a consensus approach, that was applied in high-throughput virtual screen in order to find novel $K_V1.1-(1.2)_3$ inhibitors. At first, possible binding sites for putative ligands had to be detected.

3.2.1 Blind Docking – Detection of Possible Target Sites

In order to detect putative ligand binding sites of $K_V1.1-(1.2)_3$ the *blind docking* setup as described in section 3.5 was used. The entire test-library from Xention was docked using the program Vina. Taking the fourfold symmetry of the recep-

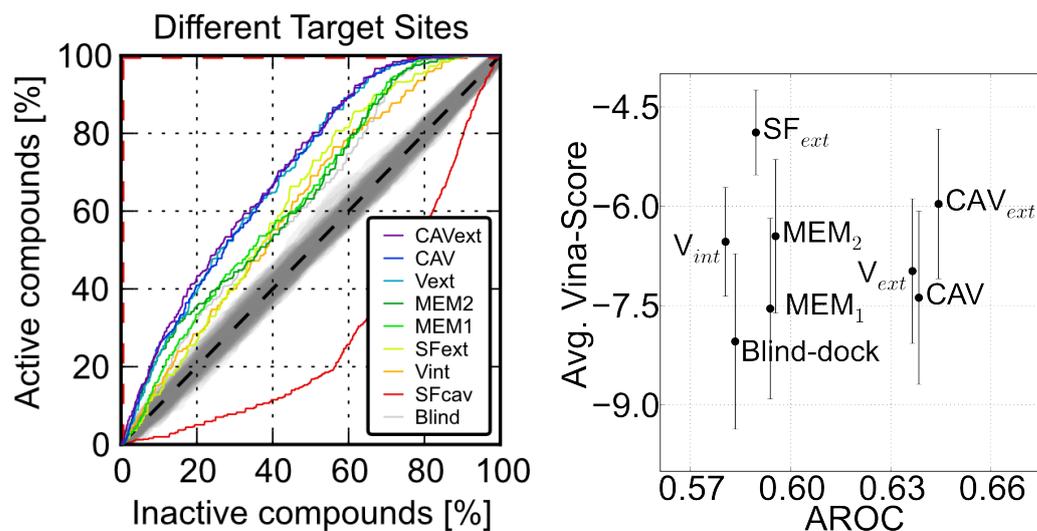


Figure 3.4: ROC-curves according to the individual sites (left). Mean scores of the Vina scores and standard deviations indicated by the error bars (right).

tor into account, blind docking established five distinct putative binding sites. Two separate sites (MEM1 and MEM2) were found at the protein-lipid surface of the TM region. Further binding sites were discovered at the intra- (Vint) and extracellular (Vext) surfaces of the voltage sensing domains as well as the inner cavity (CAV). The outer mouth of the selectivity filter was not occupied. The test-library was then docked to each of these binding sites separately. Additionally, the extracellular site of the selectivity filter was defined as a target site (SFext). Within the inner cavity, we used three different locations and sizes of active sites in order to gain further insight into possible binding modes. First, the whole cavity was defined as the receptor (CAV). Secondly, a narrow target site was defined around the inner entrance of the selectivity filter (SFcav). Finally, the region surrounding the intracellular pore entrance was defined as a binding site (CAVext).

3.2.2 Targeted Docking – Verification of Target Sites

The sites defined in blind docking were used as individual target sites. All compounds were docked into the individual binding sites, and the separation between active and inactive compounds was determined using receiver-operator-characteristic (ROC) curves. As in detail explained in section 2.3.2, the ROC curves show the fraction of identified active compounds over the fraction of inactive compounds in a ranked list ordered ascendingly. With the exception of the

Table 3.1: AROC and BEDROC values with respect to docking scores and individual scoring function terms as well as their Pearson correlation with compound mass. [a] Pearson correlation coefficients of the compound mass and docking sub-term. [b] Consensus scores of the sub-scores Vina Score, FlexX Lipo, eHiTS Strain, and Glide Evdw.

Score	PCC ^[a]	AROC	BEDROC
Rank-by-rank ^[b]	–	0.76	0.7
Rank-by-number ^[b]	–	0.75	0.69
Rank-by-max ^[b]	–	0.73	0.69
Mass	1	0.74	0.52
Vina Score	0.82	0.7	0.51
FlexX Lipo	–0.77	0.7	0.59
eHiTS Strain	–0.73	0.73	0.58
Glide Evdw	–0.67	0.74	0.68
Glide Emodel	–0.61	0.71	0.58
eHiTS Energy	–0.46	0.61	0.33

SF_{cav}, the ranks for all target sites resulted in significant enrichment of active compounds (Figure 3.4). The narrowness of site SF_{cav} led to a systematic exclusion of larger compounds, which explains the low enrichment observed. The highest enrichment was observed for the extracellular pore entry (CAV_{ext}), followed by CAV (Figure 3.4) and Vext. Among the binding sites CAV, CAV_{ext}, and Vext, the mean scores were lowest for binding site CAV. Lower scores were gained only at the transmembrane site MEM1 and for the blind docking in total (Figure 3.5). However, under physiological conditions, the compounds at MEM1 would have to compete with lipid molecules at the protein-lipid interface. Therefore, this binding site was discarded, focusing more on the binding site CAV for the subsequent steps.

The compounds from the library were further docked to the inner cavity (CAV) using eHiTS, FlexX, and Glide. Of 2675 total compounds, 2099 (including 473 active compounds) were successfully docked by all programs and were therefore included in further analysis. None of these programs enriched active compounds as significantly as Vina. The area under the ROC-curve (AROC) (Section 2.3.3) was 0.7 for Vina and less than 0.55 for each of the other approaches, represented by the FlexX-Total, eHiTS-Score, and Glide-Gscore (Table 3.1). However, as demonstrated by Stahl *et al.* [117] it is also possible to combine scoring function sub-terms, i.e. the hydrophobic term from the first algorithm and the polar in-

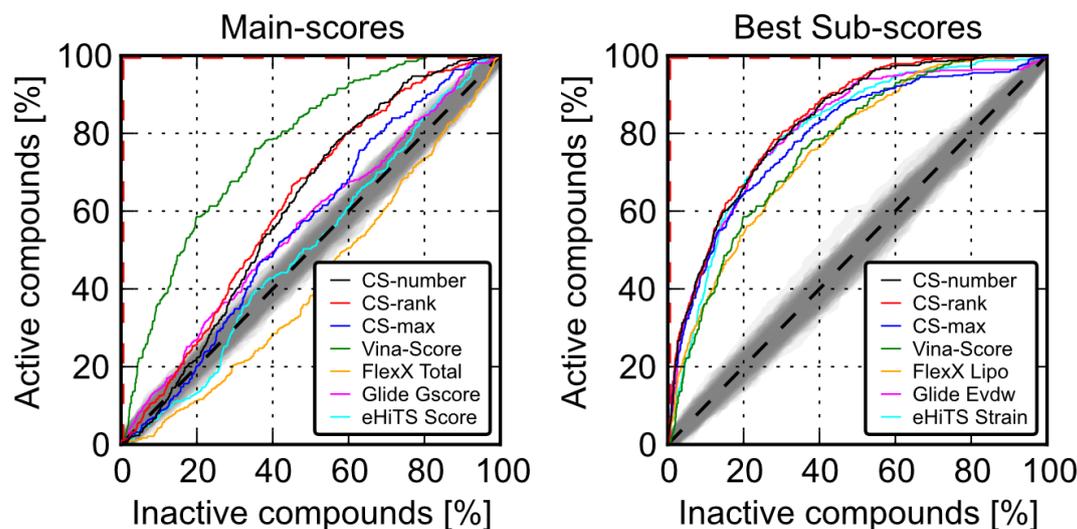


Figure 3.5: ROC-curves according to the mean scores of Vina, eHiTS, Glide and FlexX (left) and the best sub-terms of the four programs (right).

teraction term of the second algorithm. The analysis of the numerous individual sub-terms of FlexX, eHiTS and Glide revealed that some sub-terms have more predictive power than the corresponding total scores, as indicated by the ROC and BEDROC values in table 3.1. After identifying the most predictive sub-terms of the individual molecular docking algorithms, I then aimed to combine these terms in a consensus score to maximize the enrichment in the first 5% of ranked compounds.

3.2.3 Consensus Scoring – Optimization of Molecular Docking

In order to generate an optimized consensus approach, that is superior to the individual molecular docking programs, we first standardized the total scores and sub-terms using Z-scores (Section 2.3.4). Hereby, the individual scores are projected to distributions with a standard deviation of one and a mean value of zero. Three consensus scoring methods were then applied for the generation of a ranked lists:

rank-by-num: Compounds were ranked according to the mean of scores from the different scoring functions.

rank-by-rank: Compound ranks were calculated according to the individual scoring functions, then ranked according to the mean of their ranks.

rank-by-max: Compounds were ranked according to the maximum of all scores from the different scoring functions.

First, the combination of the total scores into consensus scores did not provide improved enrichment. The ROC-curves of all three CS techniques were between the ROC-curves of the individual scores (Figure 3.5). The highest enrichment from consensus scoring using the main scores of all approaches was 0.61 by rank-to-rank. However, combining individual scoring function sub-terms from the different programs, a broad range of enrichments was revealed (Supplementary table 7.1). The sub-terms leading to the highest AROC values for the individual programs were the Lipo term (0.58) for FlexX, the Evdw term (0.67) for Glide, and the eHiTS term Strain (0.58) (Table 3.1). In the case of Vina, it was not possible to check sub-terms as only the final score was provided. The use of consensus scoring methods using the Vina score in combination with the sub-terms Lipo (FlexX), Strain (eHiTS), and Evdw (Glide) led to a slightly enhanced enrichment in terms of AROC, and a significant enhancement of the BEDROC as described by Truchon *et al.*[60] (Table 3.1). The BEDROC metric is more sensitive to changes in initial enrichment as defined in the methods section. The corresponding ROC curves are shown in Figure 3.5. The consensus approaches increased the AROC value by 2 to 5% and the BEDROC value by 17 to 20%, with respect to the average AROC/BEDROC values of the individual terms used for the consensus. As indicated by the increase in BEDROC values, all three consensus approaches led to a significant enrichment in the top 8% of a ranked list of compounds. Notably, a strong dependence of the enrichment on compound mass was found, indicating a higher activity on average for larger compounds. The rankings according to the three shown consensus schemes were superior to all individual scores and sub-terms. In order to assess the quality of these consensus approach, we applied it in a high-throughput virtual screen and assayed a sample of top scored compounds.

3.2.4 High-Throughput Virtual Screening – Prediction of Novel Active Compounds

For the application of the consensus approach, introduced in the last section, two slightly different versions of the scheme *rank-by-max* were implemented and applied, resulting in two compound sets: A and B. The compounds in set A are based on prediction of a consensus score that was generated using Vina, FlexX, and Glide. The compounds in set B are additionally based on the Strain term from eHiTS and additional filtering according to drug-like prop-

erties. Initially, the *clean-drug-like* subset of drug-like compounds of the ZINC (<http://zinc.docking.org/>) chemical molecule database from 2009-11-13, containing 9,497,542 entries, was screened against the the inner cavity of K_v1.2 using FlexX. The best 20,000 compounds, according to the Lipo term from FlexX, were evaluated in the other programs as well (Glide, eHiTS, and Vina). These compounds were conducted to two different treatments:

- A) The top 20,000 structures, according to the Lipo-score from FlexX were docked with Vina and Glide. The rank-to-max consensus method was applied using FlexX's Lipo-Score, the Evdw term from Glide, and the predicted binding free energy from Vina. Only compounds commercially available from Enamine Sales (<http://www.enamine.net>) were taken into further consideration. The 200 top-ranked compounds, according to the CS scheme rank-to-max, were selected.
- B) The library of 20,000 compounds was prefiltered to remove compounds that did not fit the drug-like filter of the OpenEye FILTER software. From 20,000 molecules, 1906 were retained and screened using eHiTS. Ligand dockings were evaluated using SPROUT (version 6.3) and MAESTRO. The rank-to-max consensus method was applied using the sub-score Strain from eHiTS in addition to sub-scores from FlexX, Glide, and Vina, which were also used for the previous implementation. The top 200 compounds, according to the CS scheme rank-to-max and based on availability from Enamine Sales, were selected.

A combined list of compounds from both A and B was generated. The list was filtered to remove compounds with a logarithm of the calculated water-octanol partition coefficients (logP values) greater than 4.0, in order to ensure sufficient solubility. A total of 89 compounds were purchased from Enamine Sales. The final library of purchased compounds contained 33 ligands predicted by A and 74 compounds predicted by B; 18 compounds were common to both A and B. The final criterion for the acquisition of the 89 compounds was commercial availability.

3.2.5 Experimental Validation

In order to validate the developed consensus scoring approach, the sample of 89 top scored compounds was subjected to electrophysiological measurements, that

Table 3.2: Number and fraction of active compounds (+80% inhibition at $10\ \mu\text{M}$) from implementations A and B.[a] Fraction of active compounds in the screened subset.

	Set A	Set B	Total
Total	16829	1285	
Screened	33	74	89
Active	7	8	14
Fraction[%] ^[a]	21	11	17

were performed in the automated patch clamp setup. The cell preparation and the experimental conditions are described in detail in [128]. The measurements confirmed in total 14 compounds exhibiting $> 80\%$ inhibition of $K_V1.1-(1.2)_3$ when tested at a concentration of $10\ \mu\text{M}$. By applying this threshold we concentrated on compounds with affinities in the low micromolar range. The fractions of the identified active compounds in set A and B are 21% and 11% (Tab. 3.2). The original ranks for the 14 hits, as well as the ranks within the subset of compounds available at Enamine, are provided in the appendix (Supplementary tab. 7.3). One active compound (ID=7) was shared between both sets A and B. Assuming a uniform distribution of active compounds within the first 200 compounds of each list, the number of active compounds can be estimated to be between 11 and 39% for set A and between 3.9 and 25% for set B, with a confidence interval of 95% each. IC_{50} values for these 14 compounds lie between 0.58 and $6\ \mu\text{M}$ (Fig. 3.6). To estimate the specificity of the ligands with respect to other targets, the 14 compounds were evaluated against three important cardiac ion channels: Nav1.5, Cav1.2, and hERG. The experiments reveal a pronounced selectivity for $K_V1.1-(1.2)_3$ over the cardiac channels (Table 3.3). Notably, compounds 1 and 2 were at least 30-fold more active toward $K_V1.1-(1.2)_3$ over the other channels.

Chemical structures of the 14 active substances are shown in Figure 3.7. Physiological properties that are relevant for an estimation of their drug-like qualities are listed in Table 3.3. The drug-like scores were calculated using the Molsoft drug-likeness and molecular property estimator (<http://www.molsoft.com/mprop>). The drug-likeness model score predicts drug-like properties using Molsoft's chemical fingerprints. Values between 0 and 2 indicate very drug-like molecules, although values as low as -1 are frequently reached by drug-like molecules. Non-drug-like molecules usually give values between -3 and -0.5 . The distributions of drug-like and non-drug-like molecules are shown on the Molsoft

Table 3.3: Molecular properties and drug-like scores of the active compounds. IC₅₀ values of the 14 active compounds on K_v1.1-(1.2)₃ (primary target) and on Nav1.5, Cav1.2 and hERG (part of the cardiac safety panel). [a] Indicates the implementation that suggested the compound; Molecular weight (W); Number of hydrogen bond acceptors (A) and donors (D); calculated octanol-water partition coefficient (LogP); Molsoft’s drug-likeness model score (DL); Average concentration that produces 50% inhibition (AV) of the respective channel K_v1.1-(1.2)₃ and on Nav1.5, Cav1.2 and hERG, respective standard deviation (SD) and number of evaluations (n). [b] No SD or n are indicated for the Cav1.2 results, obtained with flexstation. Molecular properties were calculated with the molsoft molecular properties calculator.

ID	Set ^[a]	Properties and DL-score					K _v 1.1-(1.2) ₃			Nav1.5			Cav1.2 ^b		hERG	
		W	A	D	LogP	DL	AV	SD	n	AV	SD	n	AV	AV	SD	n
1	A	492	6	2	2.64	0.63	0.71	0.19	3	27.92	4.17	4	30	30	0	3
2	A	485	6	0	3.28	0.18	0.79	0.03	3	30	0	3	30	30	0	3
3	A	493	4	2	3.91	-0.67	1.41	0.24	3	24.07	10.27	3	10.44	13.26	3.97	4
4	A	482	6	1	4.34	0.42	1.62	0.48	4	28.38	2.81	3	3.3	8.05	3.53	5
5	A	427	2	2	2.76	1	2.98	0.35	4	25.73	7.4	3	30	30	0	4
6	A	492	7	2	3.56	-0.15	4.07	0.23	3	30	0	3	15.5	30	0	3
7	A,B	471	5	1	4.29	0.08	1.53	0.37	3	30	0	3	11.43	30	0	3
8	B	490	5	0	5.57	0.11	0.58	0.11	3	30	0	4	2.87	7.33	4.75	3
9	B	468	5	0	5.32	-0.32	0.93	0.48	3	29.39	1.06	3	1.55	8.62	2.74	3
10	B	495	4	0	3.89	0.37	1.66	0.21	4	30	0	4	2.52	9.46	1.23	3
11	B	477	6	0	5.59	0.27	2.71	0.64	6	22.73	9.22	4	3.85	8.43	1.88	3
12	B	480	5	0	5.09	-0.27	3.7	0.6	4	30	0	3	4.74	8.93	2.05	3
13	B	480	5	1	4.38	0.4	3.77	1.07	3	24.8	9	3	9.61	30	0	4
14	B	475	6	1	2.74	0.07	5.94	0.67	4	30	0	4	12.27	30	0	3

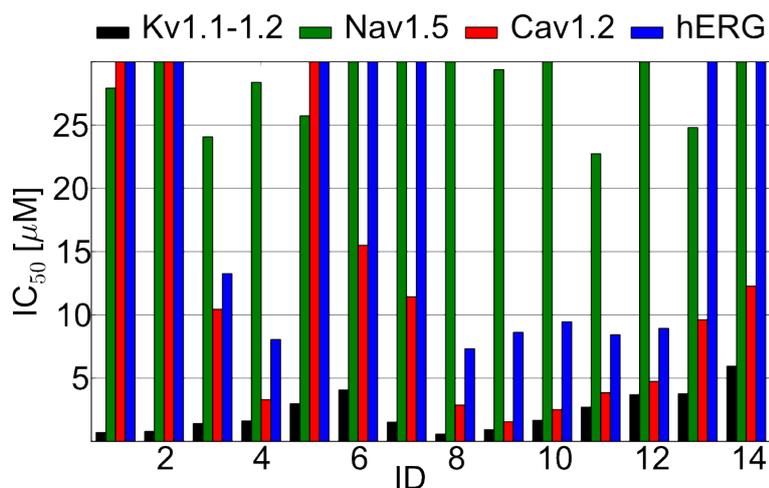


Figure 3.6: IC_{50} values of the 14 novel active compounds according to the different channels ($K_V1.1-(1.2)_3$, Nav1.5, Cav1.2 and hERG) calculated from concentration-response curve titrations.

website (<http://www.molsoft.com/mprop>). All 14 compounds share a carboxyl group close to their geometric center. Compounds 8, 9, and 12 share a Tanimoto similarity (Section 2.1.4) greater than 0.8 and have a common 4-(1,2,3,4-tetrahydroisoquinoline-2-sulfonyl)benzamide motif, which is also seen in compound 11. The similarity of compound 11 to the former compounds is 0.7 at maximum. These compounds can be regarded as one structural cluster. A second cluster comprises compounds 4, 6, 10, and 13 which each contain a 3-formylbenzene-1-sulfonamide group. Compounds 1 and 2, which are highly selective for $K_V1.1-(1.2)_3$, are not represented by either of these clusters. Twelve compounds contain a sulfur atom, and in 10 cases this takes the form of a sulfonyl group. The molecular weight of the 14 active compounds lies between 420 and 500 Da. The Tanimoto similarity between the 14 active compounds and the known active compounds from the training set was 0.56 at maximum, indicating pronounced structural diversity.

3.2.6 Receptor Flexibility

Since the value of the scoring function depends strongly on the relative placement of the interacting groups, the conformation of the target structure plays a crucial role for docking procedure and the ranking of the screened compounds. Unfavorable conformations of the target structure can hinder the docking algorithm to sample physiologically relevant poses by sterical hindrance or it can lead to an underestimation of favorable interactions. Especially, for directed polar in-

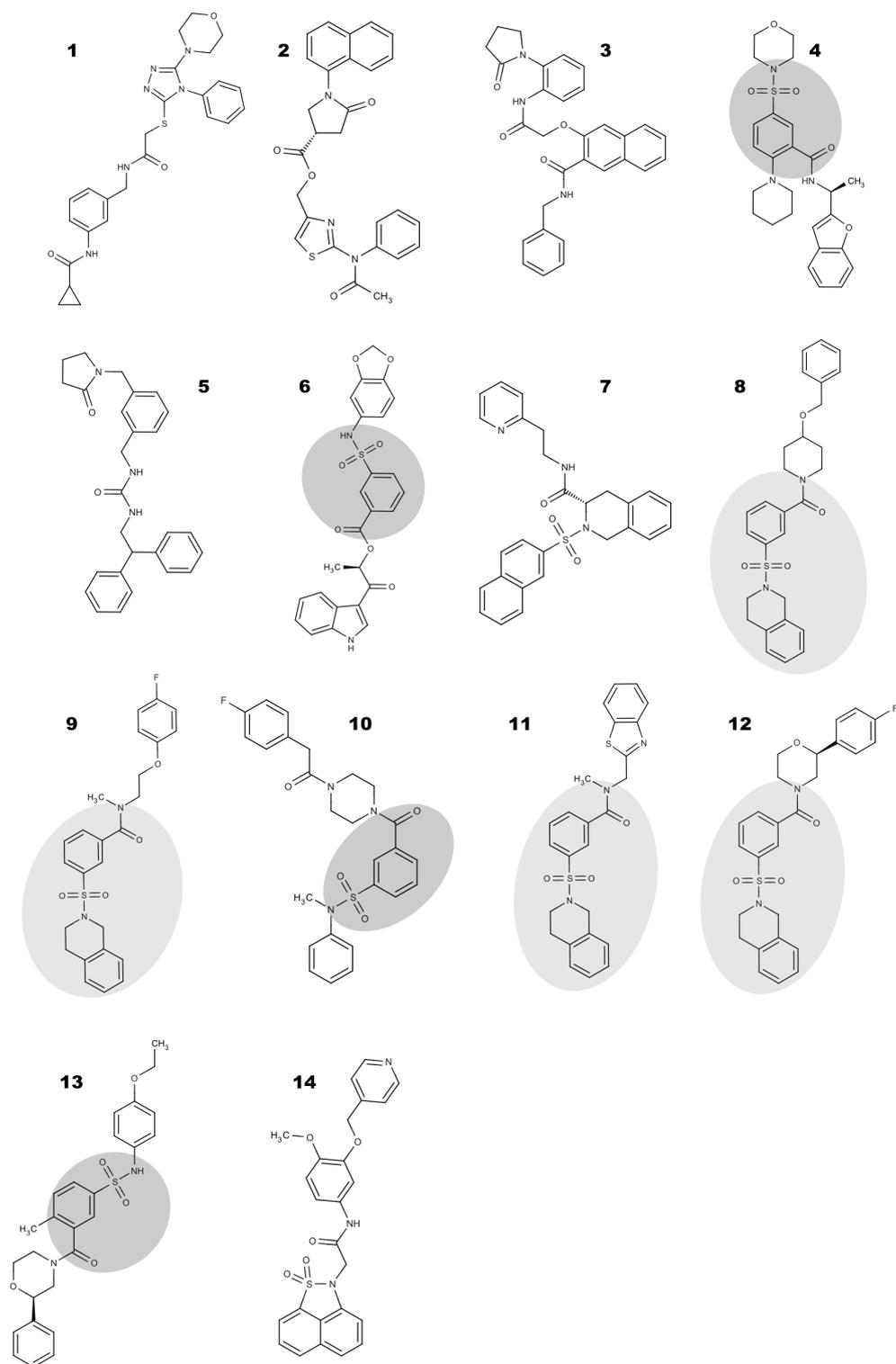


Figure 3.7: Structures of the 14 confirmed novel $K_V1.1-(1.2)_3$ active compounds. Larger repeated motifs are highlighted.

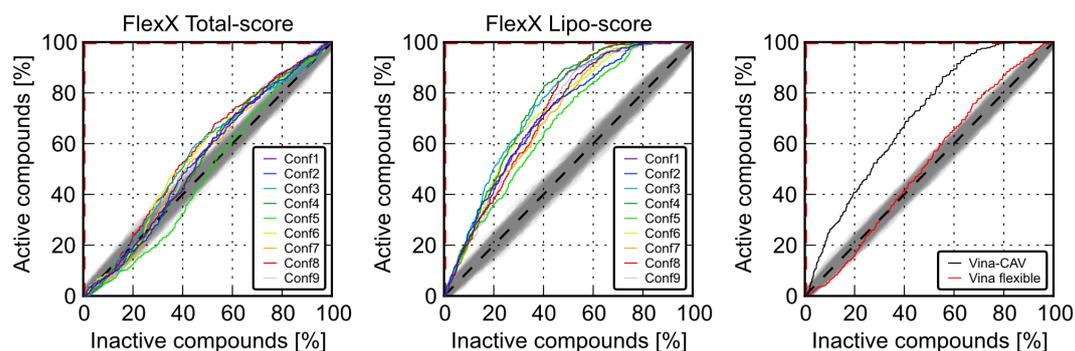


Figure 3.8: Enrichment according to different $K_V1.2$ receptor conformations using the Total score (left) and the Lipo score (middle). Docking with flexible side chains using Vina removed the enrichment (right).

teractions as H-bonds, a proper receptor conformation is important for FlexX. Therefore, other receptor conformations different from the crystal structure may exist that enhance the enrichment of our test library. A set of 10 snapshots of different conformations of $K_V1.2$ was taken from a 10 ns MD simulation trajectory. The molecular docking calculations with FlexX were repeated with these conformations serving as target structures. The ROC-curves according to the different receptor conformations are shown in Figure 3.8. As indicated by the shape of the ROC-curves, the enrichment according to the different conformations was comparable to the enrichment gained by the crystal structure.

Another way to include receptor flexibility in molecular docking is the application of approaches that allow alternative side chain (or even backbone) conformations of the receptor structure. In Vina it is possible to define the rotation of individual side chains. Therefore, the docking was repeated with Vina at the inner cavity allowing the rotation of all side chains that point into the inner cavity. Although the inclusion of side chain flexibility led to an increase of predicted binding free energy by $10 \text{ kJ}\cdot\text{mol}^{-1}$, the enrichment was completely removed in the flexible docking calculation as indicated by the ROC-curves (Figure 3.8).

3.3 Discussion

In this study, I sketched and validated a possible virtual screening protocol using molecular docking as the main technique. Four widely used molecular docking approaches have been tested for their ability to find known active inhibitors of $K_V1.1-(1.2)_3$. In this study, Autodock-Vina led to the best enrichment. Furthermore, we found that using sub-scores from the scoring functions of the individ-

ual molecular docking programs can lead to pronounced enrichments of inhibitor identification, even if no enrichment is gained using the total scoring function. Subsequent analysis indicated that the enrichment can be further enhanced by combining these sub-scores into consensus scores. These results underpin the importance of adjustment of the scoring and ranking procedures in a molecular docking calculation for successful virtual screening calculations.

The combination of blind docking with conventional docking calculations, as well as the experimental evaluation of our predictions, support the hypothesis that inhibitors bind within the inner cavity of $K_V1.1-(1.2)_3$. Using an adjusted consensus molecular docking approach, we identified several novel, potent, and selective non-peptide $K_V1.1-(1.2)_3$ inhibitors. Compounds 1 and 2 represent potential lead structures for the development of novel compounds that could selectively inhibit the ion flux mediated by $K_V1.1-(1.2)_3$ *in vivo*. Electrophysiological measurements confirmed a hit rate at or above 17% when the relatively stringent hit criteria of greater than 80% channel inhibition at 10 μ M was applied. Four compounds (1, 2, 8, and 9) bind in the sub-micromolar range (Fig. 3.7). Compounds 1 and 2 exhibit at least 30-fold greater inhibition potency toward $K_V1.1-(1.2)_3$ than they show against a small panel of cardiac selectivity targets (Nav1.5, Cav1.2, and hERG), therefore meeting some of the basic cardiac safety requirements.

Evaluation against other unrelated targets was beyond the scope of this research, but further medicinal chemistry could produce a library of similar compounds to aid in the development of a structure-activity relationship for the most active molecules, 1 and 2. This library could be used to elucidate key binding features from compounds 1 and 2 to guide the development of a novel series of $K_V1.1-(1.2)_3$ inhibitors. Neither compound 1 nor 2 has previously been reported in the literature to have any biological activity. Furthermore, there are no references associated with either of these compounds in the databases of SciFinder and Reaxys.

The high specificity, as well as the low similarity of these hit molecules to known active compounds from the training set, indicates that this approach makes proper use of the structural characteristics of $K_V1.1-(1.2)_3$ in the resulting selection process for the identification of novel structures. The drug-likeness model scores between -0.7 and 1 indicate that the 14 active compounds bear greater similarity to marketed drugs relative to non-drugs, in agreement with the fact that all compounds originate from the ZINC (<http://zinc.docking.org>) *clean-drug-like* subset. Both implementations A and B identified a similar number of inhibitors with greater than 80% inhibition at 10 μ M concentration. However, set B (74 compounds) was approximately double the size of set A (33

compounds). The fraction of active compounds was therefore nearly twice as high in A than in B. This suggests that the enrichment of active compounds is higher when the consensus scoring is applied in parallel rather than sequentially, corresponding to a situation wherein each molecular docking algorithm is applied to each compound. Although such an extensive screen would require substantially more computational time, this may prove to be the most efficient approach. Though the influence of successive filtering according to size and solubility applied after the consensus scoring procedure only in case B must be considered. Nevertheless, we show here that when the whole library was tested with only one docking program and subsequent consensus scoring was applied to a smaller library of top-ranked compounds, an improvement in enrichment of two to three orders of magnitude was achieved over a random selection of compounds [5, 115, 116, 145]. Because the consensus approach that we used in this study was trained on a library of known active and inactive compounds, this approach cannot be immediately transferred to other targets. However, it may be a reasonable starting point for ion channels that have structural and functional similarity to $K_V1.2$. Our optimization only targets the scoring and the ranking stages of molecular docking and does not affect the sampling stage. Further improvement might be possible when all three stages are included in the training process.

Notably, the inclusion of receptor flexibility did not lead to an improved prediction. The enrichment regarding alternative receptor structures in combination with FlexX resulted in comparable ROC-curves. The calculation with Vina using flexible side chains completely removed the enrichment (Figure ??). A possible explanation for these observations could be the fact that the flexibility of the side chains in the case of FlexX changed the overall conformation only slightly, while, in the case of Vina, the flexibility led to receptor conformations that were not observed in the simulation. For example, the side chains of threonine 375 were turned by 180 degrees (data not shown). In this conformation the carbonyl groups of the threonine side chains pointed to the center of the pore. Under physiological conditions one or two potassium ions are expected in the selectivity filter. These ions were not included in the molecular docking calculations. The interactions with the ions and the solvent molecules are thought to stabilize the conformation seen in the crystal structure by polar interactions. Since the ions and solvent atoms were absent in the molecular docking calculation the selectivity filter was not stabilized by these groups, allowing the docking algorithm to relax the structure to an unphysiological conformation. In contrast, ions and solvent were present in the MD simulation. Accordingly, the enrichment by the

calculated scores was significantly reduced in the case of Vina but not in the case of FlexX, where flexibility was introduced by docking to conformations that were generated by free all-atom MD simulations.

3.4 Summary

In this collaborative study I established novel inhibitors of the potassium channel $K_V1.1-(1.2)_3$. The combined use of four virtual screening (VS) programs (eHiTS, FlexX, Glide, and Autodock-Vina) led to the identification of several compounds as potential inhibitors of the $K_V1.1-(1.2)_3$ channel. From 89 electrophysiologically evaluated compounds, 14 novel compounds were found to inhibit the current carried by $K_V1.1-(1.2)_3$ channels by more than 80% at $10 \mu\text{M}$. Accordingly, the IC_{50} values calculated from concentration-response curve titrations ranged from 0.6 to $6 \mu\text{M}$. Two of these compounds exhibited at least 30-fold higher potency in inhibition of $K_V1.1-(1.2)_3$ over a set of cardiac ion channels (hERG, Nav1.5, and Cav1.2), resulting in a profile of selectivity and cardiac safety. Notably, the assayed compounds were purely selected on the basis of computational results, driving this approach fully automatable. Therefore, the results presented herein provide a promising basis for the development of novel selective ion channel inhibitors, with a dramatically lower demand in terms of experimental time, effort, and cost than a sole high-throughput screening approach of large compound libraries.

3.5 Methods

Tanimoto-coefficients section were calculated using cheminformatics and machine learning software RDkit (<http://www.rdkit.org>) and default 2048 bit hash Daylight topological fingerprints (Section 2.1.4). The minimum path size was 1 bond, the maximum 7 bonds. The quality of the predictions was evaluated by comparing the predictions of the individual programs with the experimental data in the benchmark library. The predictions were illustrated using ROC-curves and were quantified using the boltzmann enhanced discrimination of the ROC-curve (BEDROC) as well as the area under the ROC-curve (AROC) (Section 2.3.3). A weighting factor of $\alpha = 20$ was used for all evaluations, corresponding to 80% of the score from the top 8% of the list. Both the AROC and the BEDROC metric provide values between 0 and 1. The curves in figure 3.5 correspond to the same set of compounds.

The receptor input files for FlexX [101] (version 3.1.4) were generated using an in-house Python (<http://www.python.org>) script defining all atoms of the inner cavity within a cylinder of 10.5 Å around the fourfold symmetry axis as the active site. High-throughput screening was performed using FlexX. Standard parameters were used for weights of the scoring function and the number of intermediate solutions for each fragment. Autodock-Vina is in detail described in section 2.2.3. Input files were generated using the AutoDock plug-in [113] for PyMOL [28]. For blind docking, a cubic box containing the complete Kv1.2 transmembrane domain was used. Ligand clusters were defined manually by visual inspection. For targeted docking, rectangular boxes with edge lengths between 1.2 and 3.4 nm around the center of the individual ligand clusters were used. High-throughput docking was performed using 20,000 compounds from the FlexX calculation with the inner cavity as the target site.

The Electronic High-Throughput Screening programme (eHiTS) [146, 147] calculates a score for each structure according to the poses of the ligand and the complementarities of *surface points* on the ligand and the receptor. In the first instance, the binding pocket is defined by a steric grid which divides regions into separate pockets and all possible interaction sites are identified. In the next step, the ligand is divided into rigid fragments and connecting flexible chains. Then, each fragment is docked to every possible place in the binding site. Ring systems are considered rigid and their conformation is preserved as given at the input, usually the lowest energy conformer. For the recognition of interesting poses a simple and fast chemical flag based statistical scoring function is used. Afterwards, an exhaustive matching of compatible rigid fragment pose sets is

performed to yield complete structures. Typically, the program evaluates several million mappings of the rigid fragments to the target site. The reconstructed solutions define a rough binding pose which is then refined by a local energy minimization. The resulting poses are then evaluated by a third, more time consuming scoring function that combines both statistical and empirical components, plus additional grid based geometrical terms as well as entropy loss estimation. Additional terms are used in the final scoring function to further reflect all factors involved in binding, such as steric clashes, depth of the cavity, solvation, conformational strain energy of the ligand, intramolecular interactions in the ligand, and entropy loss due to hindered rotatable bonds. This final scoring function attempts to estimate the free energy of binding. The result of the final elaborated scoring function is used to rank the generated solutions.

Glide 5.5 [41, 50] performs a gradual guided progression solution space search by an initial rough estimate of the ligand conformation and a subsequent torsionally flexible energy optimization on a non-bonded potential grid based on the OPLS-AA force field [68]. The best candidates, as defined by the scoring function, are further refined by Monte Carlo sampling of the ligand pose. Glide's scoring function is a combination of empirical and force-field-based terms. Intermolecular interactions were pre-calculated on a grid representing the extracellular half of the receptor and were centered on selected residues in the binding site in such a way as to enable access to total available space in the inner cavity and include long range interactions up to 20 Å. Receptor flexibility was derived by in place temporary alanine mutations and van der Waals (vdW) radii scaling. The 20,000 ligands selected with FlexX screening were docked with full flexibility on the grid. For each ligand, ten poses were generated and subsequently clustered (RMSD < 0.5 Å).

Sequence and flexibility in MD simulation. The crystal structure of K_v1.2 (PDB code: 2A79) from rat served as template for modeling of the target structure [85]. The part of the sequence that was part of the model is shown in figure 3.9. Loops that were not present in the crystal structure were added using the MODELLER software [37]. For Glide, the protein structure was optimized with MacroModel (OPLS2005 force field), and the Protein Preparation Wizard was used to optimize hydrogen bonding networks of the protein. For FlexX the structure was energy minimized in the Amber ff99SB forcefield [57].

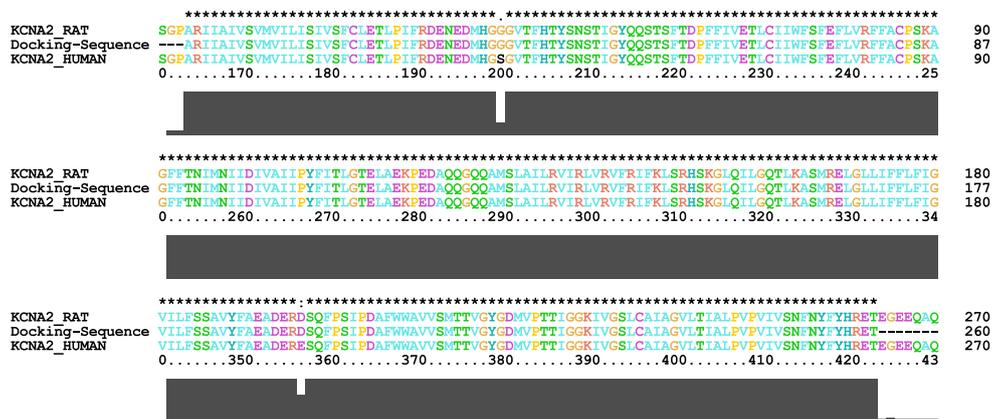


Figure 3.9: Cutout of the sequences of the rat and human $K_V1.2$ and of model that was used for the molecular docking computations.

4

Inhibition of human Aquaporin 9

In this chapter, I present a collaborative study of small molecule interactions with the human water channel protein Aquaporin 9 (hAQP9). This protein is a member of the Aquaglyceroporin family. The computational approaches that I used here, range from homology modeling to molecular docking to molecular dynamics simulations and map a complete pathway from an unknown protein structure to a detailed structural model of receptor-ligand interactions. An essential element of this study lies in the cross-fertilizing character of computational and experimental findings that guided and supported each other. All reported binding assays and the generation of the protein mutants were performed by Dr. Michael Rützler at Aarhus University (Denmark).

4.1 Introduction

In all living cells, regulation of solute and water movement across cell membranes is of critical importance for osmotic balance. Aquaporins are membrane channel proteins which facilitate the permeation of water and small neutral solutes across biological lipid membranes. Several high-resolution X-ray structures of AQPs revealed a conserved homo-tetrameric structure where each subunit provides an independent pore for water and other solutes [6, 45]. In contrast to tetrameric potassium channels, each monomer constitutes an independent pore.

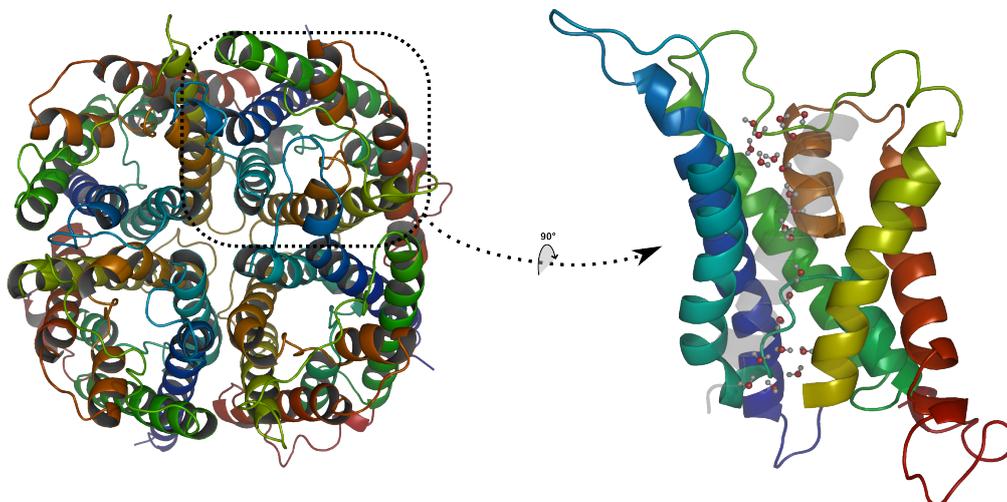


Figure 4.1: Ribbon representation of AQP1 tetramere (left) and a single subunit (right). Transmembrane domain TM3 is transparent to allow the view to the single file water layer in the pore. The conformation was taken from a equilibrium all-atom molecular dynamics simulation.

Each monomer consists of six transmembrane (TM) domains (TM1–6) connected by five loops (A–E) with intracellular N- and C-termini. The loops and the transmembrane domains form a two-fold repeated tandem structure that is referred to as hourglass structure (Figure 4.1). The first and second part of the protein share considerable sequence homology resulting in this quasi two-fold structural symmetry in the plane of the membrane [129]. Loops B and E fold into the bilayer from opposite sides flanked by the TM domains. Both loops touch each other in the middle of the pore. They contain a helical part and the Asn-Pro-Ala (NPA) signature motifs, which are highly conserved among the AQP family [1].

Both Asn residues from the NPA motifs are the capping amino acids at the end of the α -helices and form a constriction site, called the NPA-region. Close to the NPA-region at the exoplasmic side, aquaporins have a highly conserved region that consists of a ring of aromatic side chains and a highly conserved Arg residue. This so called aromatic/arginine (ar/R) region forms the narrowest part of the pore and serves as a selectivity filter for neutral solutes [6]. It has been demonstrated for AQP1 that mutation of the corresponding Arg residue (Arg195) to Val conveys proton leakage [6].

Human Aquaporin 9 (hAQP9) is a member of the Aquaglyceroporin family and has a broad solute permeability including glycerol, polyols, carbamides, purines,

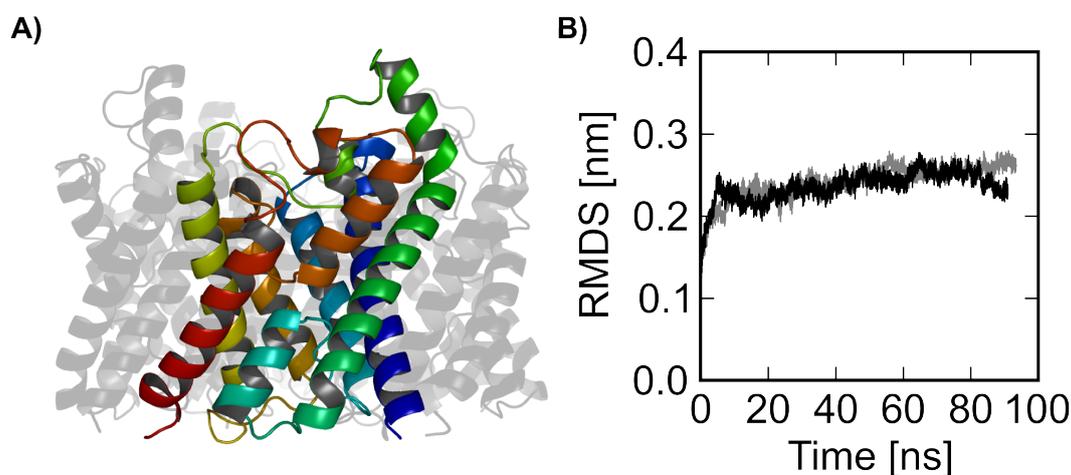


Figure 4.2: Homology model of human AQP9. Ribbon graphical illustration of the homology model of hAQP9. Only one subunit is highlighted for clarity (A). The RMSD trace in two independent MD simulations of hAQP9 indicates structural stability of the model in the simulated time scale (tens of nanoseconds) (B).

pyrimidines, nucleosides, and monocarboxylates and urea [122]. It is expressed in the plasma membrane of hepatocytes in the liver, which is a major site of production and elimination of metabolites such as urea, nucleotides, and ketone bodies. Substantial amounts of these solutes must rapidly cross the hepatocyte plasma membrane with minimal osmotic perturbation. Recently, Jelen *et al.* [62] reported the identification of novel murine AQP9 inhibitors by a small molecule screen, using an fluorescence intensity assay [40, 73, 95]. In this assay, fluorescence intensity of the calcein fluorophore is affected by cell volume changes. Consequently, with knowledge about the relation between extracellular osmolarity, fluorescence intensity and cell volume, cell volume changes can be monitored as changes in fluorescence intensity [40]. hAQP9 and the bacterial glycerol facilitators (GlpF) share a sequence identity of 33% enabling the construction of a homology model of hAQP9 using the high resolution structure of GlpF (PDB:1FX8) as template [42]. Using this structural model, I aimed to identify novel inhibitors that can serve as selective hAQP9 inhibitors *in vivo* or even as lead compounds.

In this study, a homology model of the human form of AQP9 was constructed and used for a structure based virtual screening (SBVS). Novel inhibitors of hAQP9 were identified whose inhibition activity was confirmed in the described fluorescence intensity biological assay system. Furthermore, the putative molecular

interaction site in AQP9 was identified using a combination of *in silico* screening, mutagenesis and molecular dynamics simulations. The results of this study are presented and discussed in the following paragraphs.

4.2 Results

4.2.1 Homology Model of human Aquaporin 9

Since the structure of AQP9 has not yet been resolved, I built a homology model based on the crystal structure of the bacterial glycerol facilitator GlpF (Protein Data Bank code 1FX8) as described in [62]. This model was subjected to two 100 ns atomistic molecular dynamics simulations. The simulation setup is described in detail in the Methods section 4.5 of this work. During the first 10 ns the protein atoms were position restrained, allowing both lipid and solvent molecules to equilibrate. Afterwards, the system was simulated without position restraints. During the unrestrained simulation, the root mean square deviation (RMSD) of the hAQP9 backbone with respect to the initial structure converged to a value of 0.3 nm (Figure 4.2).

4.2.2 Single Pore Water Permeability Coefficients

The mouse and human variants of AQP9 protein share a sequence identity of 75.6%. Therefore, I hypothesized that binding sites for some previously identified murine (m)AQP9 inhibitors [62] can be detected on hAQP9 as well. Putative binding sites were identified by screening the whole hAQP9 homology model in a blind molecular docking calculation. This screen revealed several distinct putative binding sites (Figure 4.3). All four pores of the tetramer were occupied by compounds from both the extra- and the intracellular side. Compounds were also placed into the inner cavity in the center of the tetramer as well as at the protein-membrane interface. I considered the pore entries as most plausible active sites for inhibitors and used the docking poses as starting configurations in MD simulations. Since the assay probes activity (water permeability) rather than affinity, I calculated the single channel osmotic water permeability coefficients (p_f) to quantify the inhibitory effect of the compounds. Simulations were carried out under physiological temperature and pressure conditions and with a KCl concentration of 124 μ M. The protein was embedded in a dimyristoylphosphatidylcholine (DMPC) lipid bilayer mimicing the natural environment of hAQP9. I performed two individual sets of simulations with compounds located either at the extracellular or the intracellular pore vestibule and calculated

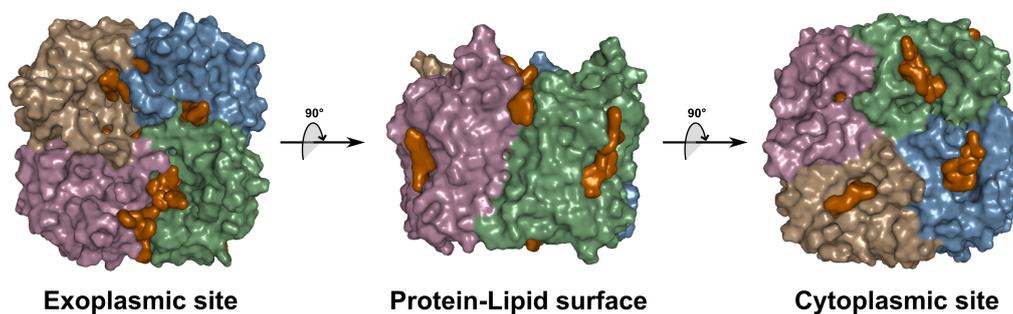


Figure 4.3: Illustration of the blind docking results. Lipid molecules (orange) were docked to the surface of the protein and revealed several possible binding pockets.

the p_f values by the mobility of the water molecules in the pore. I found that the majority of compounds led to a significant reduction of computed p_f values when placed at the extracellular AQP9 site (Figure 4.4). Compounds located at the extracellular site suppressed water permeability more (average $p_f = 0.96 \pm 0.55 \times 10^{-14} \text{ cm}^3\text{s}^{-1}$) than binding to the intracellular site (average $p_f = 1.5 \pm 0.85 \times 10^{-14} \text{ cm}^3\text{s}^{-1}$).

However, the computed p_f values of the individual molecules did not correlate with measured activities on mAQP9 (data not shown). Therefore, I decided to utilize the simulation trajectories to identify the residues in AQP9 putatively involved in the binding of the compounds. For this purpose, I monitored contacts between individual amino acids and the bound-ligand every 50 ps and summed up the counts of all simulations. The analysis revealed 6 extracellular and 4 intracellular residues which were significantly more frequently involved in contacts with the ligands than other residues (Figure 4.5). Based on this analysis we concluded that Ile60, Tyr151 and Leu209 might be involved in AQP9-inhibitor binding at the extracellular side of the channel. Furthermore, His82, Met91 and Phe180 might be involved in AQP9-inhibitor binding at an intracellular site. All numbers relate to the human AQP9 primary protein sequence.

Experimental Validation

In order to distinguish between these two alternative sites, I suggested Ile60, Tyr151 and Leu209 at the extracellular site and His82, Met91 and Phe180 at the intracellular site for mutagenesis. According to this suggestion, conservative changes were introduced, based on homologous AQP sequences, into hAQP9 by site-directed mutagenesis. The changes made are summarized in Tab. 4.1.

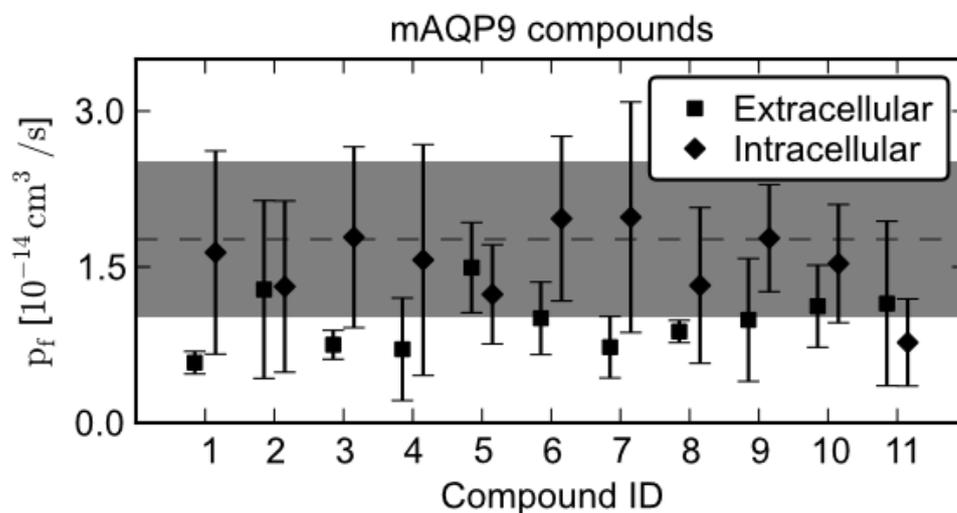


Figure 4.4: Computed p_f values from MD simulations with the hAQP9 homology model in complex with mAQP9 inhibitors. Compounds were ranked according to their inhibitory effect of mAQP9. The values represent the mean from four individual subunits of the tetramer and the error bars indicate the standard deviation. As a reference, the p_f computed from two separate MD simulations without inhibitors is depicted with a dashed line (average) and the standard deviation indicated by the grey background.

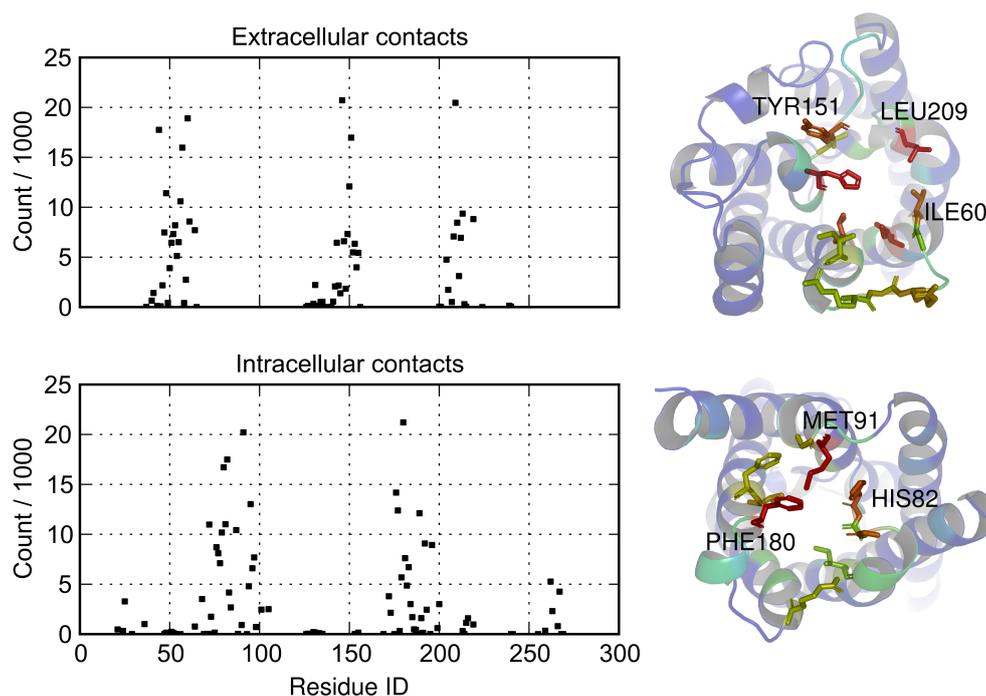


Figure 4.5: Frequency distribution of simulated contacts between inhibitors and hAQP9 amino acid residues. Contacts within the 150 ns time window of simulation between inhibitors and hAQP9 were analyzed and are displayed along the primary amino acid sequence as contacts per 1000 total contacts made (left). The illustration of the putative inhibitor interacting amino acids in the homology model viewed from the extracellular side (top) and the intracellular side (bottom). Residues in the illustrations are colored according to the observed number of contacts from blue (no contacts) to red (highly frequented).

Table 4.1: Conservative amino acid changes introduced into hAQP9 by site-directed mutagenesis.

Mutant	Site	Seq. in hAQP9	Template AQP	Homologous seq.
I60V	extracellular	VITINV	mouse AQP7	YLGVNL
H82A	intracellular	SGGHIN	yeast AQPY1/2	SGGALN
M91N	intracellular	SLAMCL	mouse AQP7	TFTNCA
Y151F	extracellular	FATYPA	no template	
F180V	intracellular	FAIFDS	mouse AQP3	LAIIVDP
L209M	extracellular	SLGLNS	mouse AQP7	SLGMNS

Subsequently, three individual stable chinese hamster overy (CHO) cell lines were established for each construct. In a first set of experiments the functionality of these constructs in osmotic water permeability assays was determined. These experiments indicated that 5 out of 6 hAQP9 mutant genes encoded fully functional proteins in this assay. Cell lines expressing hAQP9 F180V displayed reduced water permeability, compared to hAQP9 expressing cells, when expression was fully induced (by tetracycline). This suggested that hAQP9 F180V might be affected in single channel water permeability, protein stability or protein trafficking. Therefore, these cell line was excluded from further analysis. In a next step it was tested if any of the mutations alter the inhibitory effect of the studied ligands. It was found that the two amino acid exchanges made at the intracellular side altered the inhibition, whereas the exchanges made at the extracellular side did not (Figure 4.6). I therefore concluded that several identified AQP9 inhibitors, including phloretin, CD05595, RF03176 and HTS13772 likely interact with AQP9 in close proximity to His82 and Met91, both located at the intracellular side.

4.2.3 High-Throughput Virtual Screening

The mutational analysis indicated that the binding site for the studied ligands is located at the intracellular entrance of the pore of hAQP9. I therefore concentrated on the intracellular binding site for subsequent docking studies. To probe if this site can be used to predict the experimental observations when used as target site in a molecular docking calculation, I generated a test-library of compounds with 11 true hAQP9 active compounds and around 2400 random de-

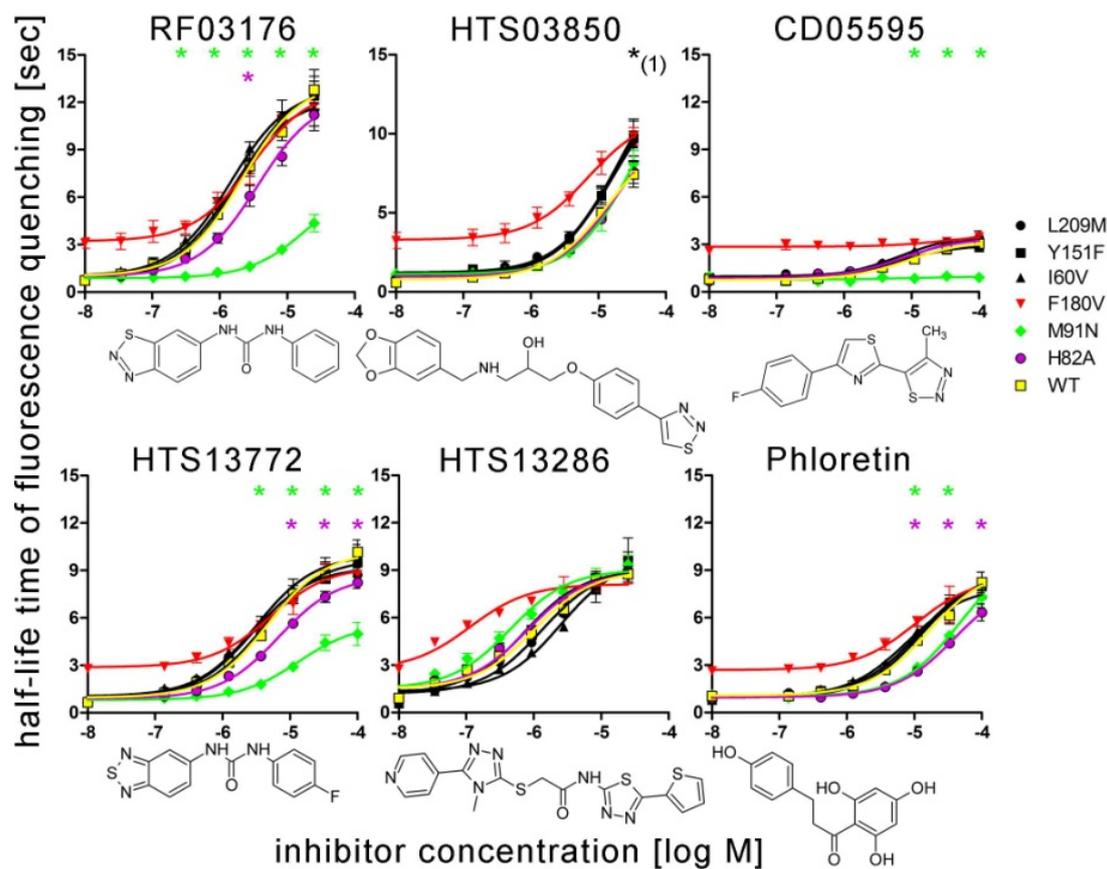


Figure 4.6: Effects of hAQP9 single amino acid mutations on inhibitor potency in CHO-hAQP9 cell shrinking assays. Induction of hAQP9 was titrated with tetracycline to achieve similar baseline cell water permeabilities in all cell lines, except for CHO-hAQP9 F180V, where this was not possible. The lowest x-axis values in each dose-response curve are 0, and have been altered for presentation on a log-transformed graph.

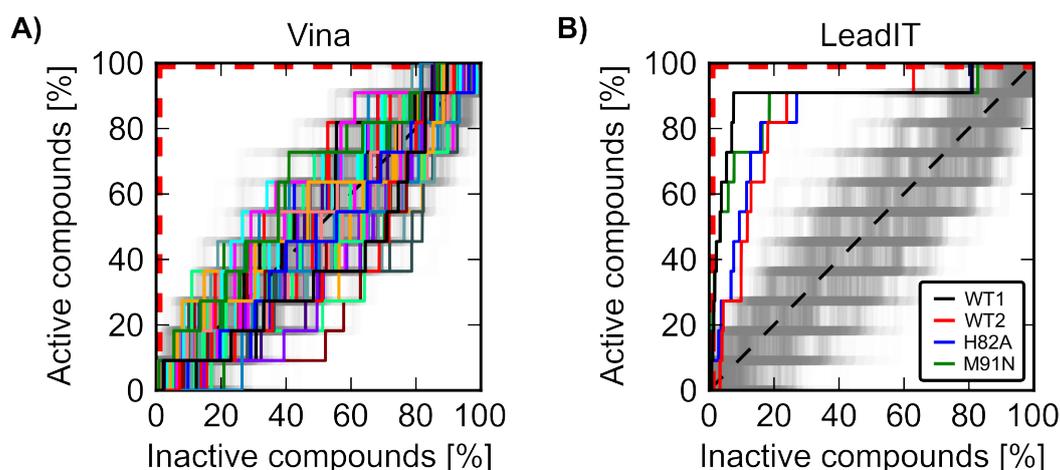


Figure 4.7: Benchmark results ROC-curves for Vina using rankings according to multiple receptor configurations (left) and for LeadIT using two different conformations of the WT and single structures of the mutants H82A and M91N (right). The steps indicate the ranking of the 11 active compounds in the test library.

coys. The content of this test-library was docked using the program Vina against each of the receptor configurations. The ROC-curves generated from the rankings according to the individual calculations lie in the random regime without exception (Fig. 4.7). Therefore, no enrichment of active ligands was observed. Next, I changed the molecular docking algorithm and used the program LeadIT. Due to a restricted number of licenses that would allow the parallel calculation against multiple targets, I only used two different configurations of the hAQP9 homology model. The first was the energy minimized homology model (WT1) and the second a configuration from a MD simulation (WT2). The compound rankings according to the docking calculations were analysed by generating the respective ROC-curves (compare section 2.3.2). Both curves significantly sample the non-random region. Therefore, LeadIT achieved significant enrichment for both of the receptor configurations of the wild-type hAQP9. The molecular docking calculation was repeated with the mutants H82A and M91N. The enrichment according to the mutant structures was comparable to the enrichment on the hAQP9 wildtype models (WT1, WT2). Accordingly, I decided to use the program LeadIT in a high-throughput virtual screen for novel hAQP9 inhibitors.

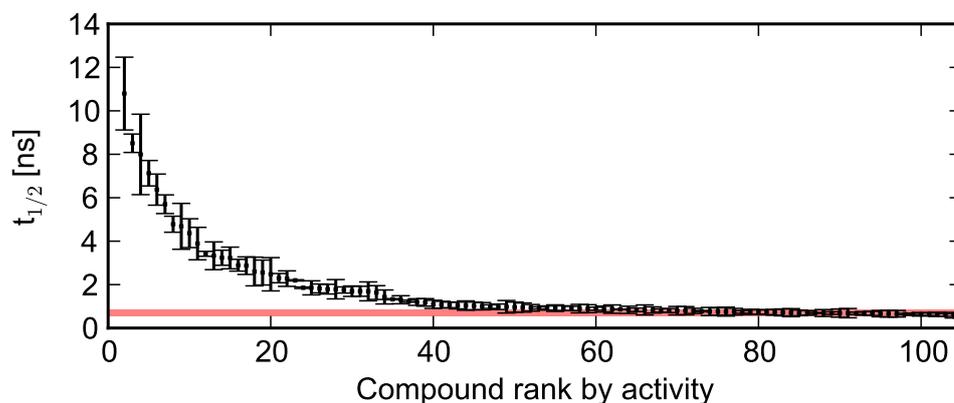


Figure 4.8: Obtained shrinking half-times ($n=3$) of the 105 tested compounds, and standard deviations (error bars). The red line indicates the value according to pure DMSO. Higher shrinking times indicate higher inhibitory efficiency.

4.2.4 Identification of Novel Inhibitors

According to the results of the benchmark presented in the previous section, I used LeadIT and the model of the hAQP9 wild type to screen around 1 million compounds *in silico*. The 105 top ranked compounds of this screen were purchased and activities were measured in the cell based assay. The effects of these compounds were tested on CHO-hAQP9 cell water permeability. We found that 30 out of 105 tested substances conferred significantly reduced water permeability in CHO-hAQP9 cells, compared to DMSO treatment (Fig. 4.8 and supplementary Table 7.6). Furthermore, by performing dose-response analyses for the top 18 ranked (by activity) compounds, we found that the most potent inhibitors affected CHO-hAQP9 cells at an apparent half-maximal (IC_{50}) concentration between 4 to 10 μM (Supplementary Table 7.8). To further validate the predictive AQP9-inhibitor interaction model, dose response analyses were performed for the 6 best (by IC_{50}) compounds. Here we found that 5 out of 6 substances are affected by at least one of the amino acid exchanges H82A or M91N, suggesting inhibitor-hAQP9 binding near these amino acid residues (Figure 4.9), and therefore the target binding site.

Computational Analysis

The predicted binding poses of the top 6 active compounds are complementary to the chemical environment in the pore and span the NPA region (Figure 4.10). Some compounds range to the extracellular site and show interactions

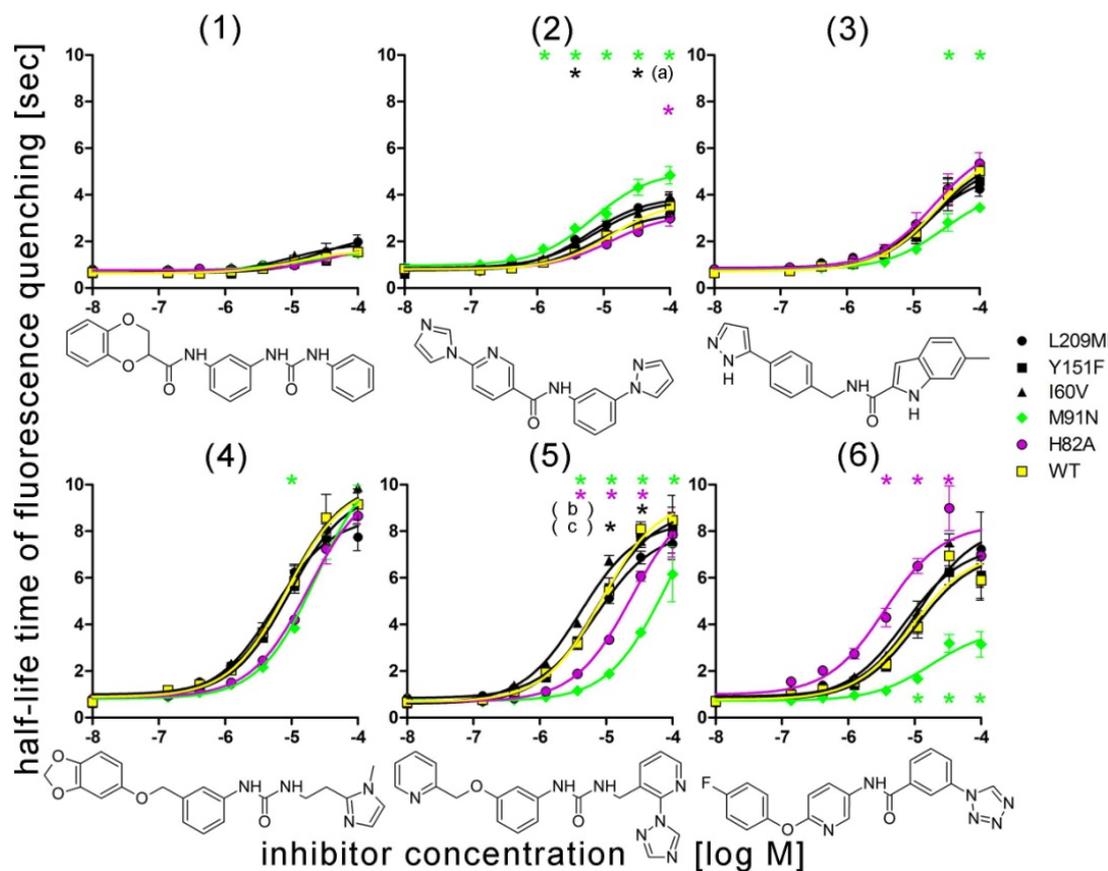


Figure 4.9: Effects of hAQP9 single point mutations on CHO-hAQP9 cell water permeability. Induction of hAQP9 was titrated with tetracycline to normalize the cell water permeabilities in the cell lines. Intracellular amino acid exchanges resulted in enhanced or reduced inhibitor potency for 5 out of 6 tested substances, corresponding to the 6 substances with the lowest apparent IC₅₀ (see supplementary spreadsheet 2). Of these, compound (1) was a very moderate inhibitor of CHO-hAQP9 cell water permeability.

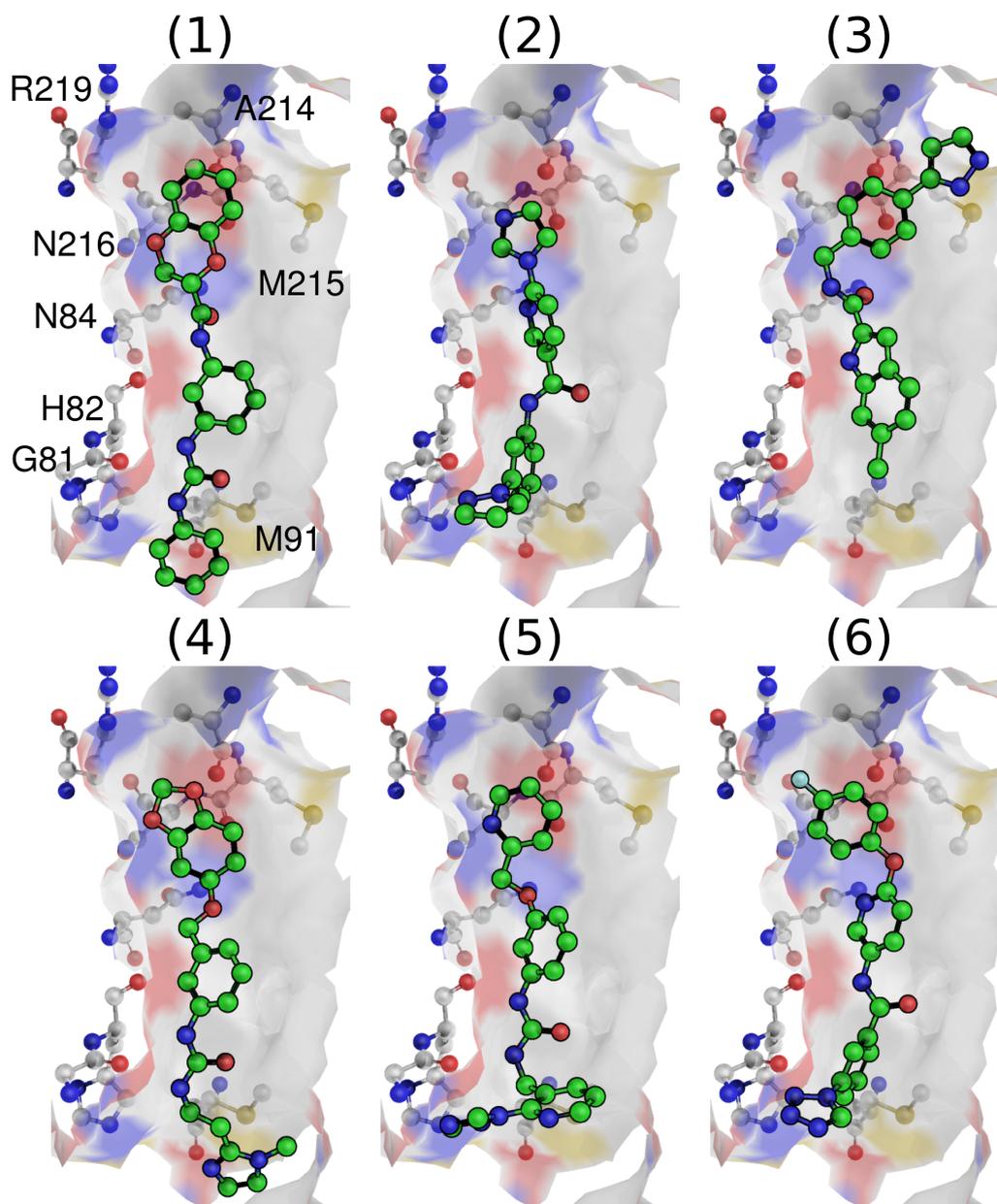


Figure 4.10: Putative positions of the top 6 tested substances (ID1–6) (in ball-and-stick representation) docked to the intracellular side of the pore of the hAQP9 homology model represented as ball-and-sticks as well as surface. The colors decode the chemical elements, oxygen (red), nitrogen (blue), sulfur (yellow). For clarity, the carbon atoms of the ligands are colored green, whereas the carbon atoms of the receptor are colored grey. The ID's correspond to the ID's in figure 4.9. The compounds are ranked according to the calculated IC_{50} values.

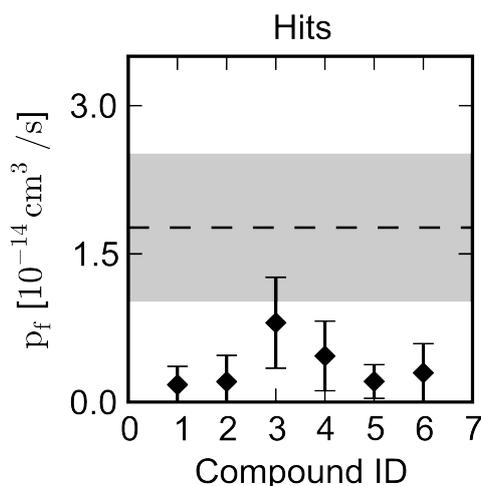


Figure 4.11: Calculated p_f values from MD simulations using the docking poses as initial configurations for the substances (ID1–6) from Fig. 4.9 and 4.10. The values represent the mean from four individual subunits of the tetramer and the error bars indicate the standard deviation. As a reference, the p_f computed from two separate MD simulations without inhibitors is depicted with a dashed line (average) and the standard deviation indicated by the grey area.

with Arg219 from the ar/R region. These six compounds are predicted to be in very similar poses in contact with the side chains of His82 and Met91. An urea group, which is a regular motif of these compounds, was observed interacting with the backbone carbonyl groups of Gly81 or His82. Furthermore, all poses predict the presence of an aromatic group between the NPA motifs. Moreover, the bound molecules show a hydrogen acceptor in close vicinity to Asn84 of the intracellular NPA motif as well as a donating group in vicinity of Gly81 or His82. I simulated each of these compounds for 30 ns and used the last 20 ns to calculate p_f values. Notably, in all the cases where the compounds were present, the calculated p_f values were significantly smaller than the reference value without compounds (Figure 4.11) as well as the p_f values from the previously simulated set (Figure 4.4), therefore confirming significant inhibition.

4.2.5 Simulated Ligand Association

Although AQP9 is a member of the aquaglyceroporin sub-family, whose members have wider pores than the pure water channels, it was questionable whether the pore is wide enough to allow the entry of drug-sized compounds. In contrast

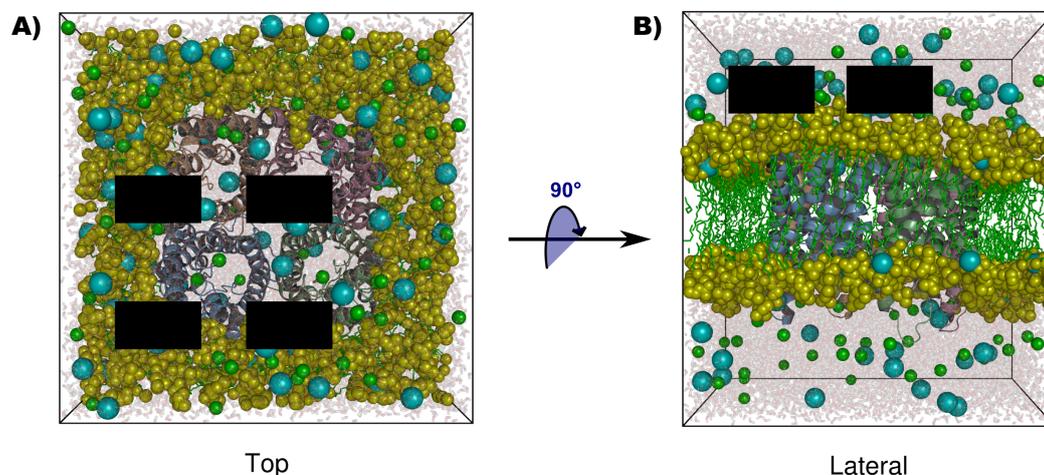


Figure 4.12: Simulation box with the solvated model of AQP9 (cartoon representation) embedded in a DMPC lipid bilayer (green and yellow). Potassium and chloride ions are represented as spheres (green, cyan). The molecules CMP1 are covered by black boxes.

to the docking approach described in section 4.2.2, where the compounds were docked rather superficially to the protein surface (data not shown), the current compounds are predicted to bind deeply into the water pore. To address the question of accessibility of the predicted poses, I set up a set of MD simulations. The goal of these simulations was also to simulate the spontaneous occupation of the intracellular pore entrance and to identify possible binding pathways as well as possible other modes of binding.

From all compounds in our library that were affected by AQP9 single amino acid exchanges, we selected a compound with a relatively low half maximal inhibitor concentration ($IC_{50}=0.39 \mu\text{M}$) and a relatively high maximal inhibition rate (90%) for human AQP9. I refer to this compound as CMP1. Due to a currently ongoing patent application, the structure can not be shown in this work. The simulation box was set up containing 4 ligand molecules that were placed in the bulk solution at least 1.5 nm away from the hAQP9 tetramer (Figure 4.12). 40 simulations were carried out with different initial velocity distributions. In 30 simulations one or more inhibitor molecules diffused into the membrane bilayer. In two cases, a spontaneous binding at the intracellular site was observed. The configuration generated by the MD simulation resembles the predicted pose from the molecular docking calculation within 0.5 Å RMSD. In these simulations, a

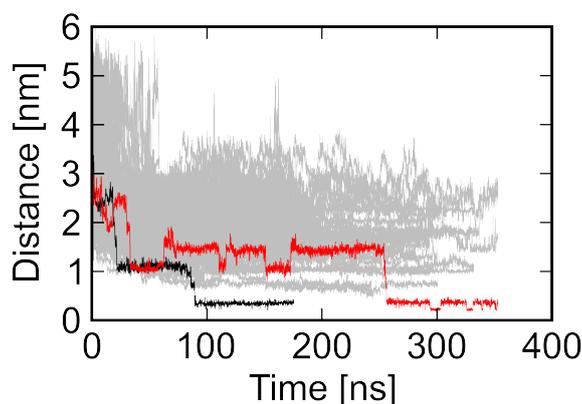


Figure 4.13: Minimal distance of compound CMP1 to the intracellular NPA motif. Two trajectories (black and red) spontaneously converge to the pose predicted by molecular docking.

single compound entered the pore of the hAQP9 homology model at the intracellular site approaching the Arg219 of the ar/R region by less than 5 Å. Notably, CMP1 reached a configuration and orientation close to the predicted poses of the 6 hits (ID1–6) presented in the last section, and spanned the NPA region of hAQP9. During the process of binding the molecules entered multiple metastable transition states (4.13). The entering demonstrated that the pore is large enough to allow the entering of drug-sized molecules. In both simulations where a binding was observed, the actual entering of the pore, from the outer mouth to the binding site, happened within less than 10 ns simulation time (4.13). Also ligand configurations were observed which were significantly different from the molecular docking prediction. Notably, also alternative binding modes were generated by the simulations. In addition, also extracellular associations of CMP1 to the hAQP9 homology model were observed. In two cases, one compound occupied the center of the tetramer from the extracellular site. However, the association of single CMP1 molecules at the extracellular site (neither at the pore entry nor at the central cavity) did not lead to a complete inhibition of water flux through the pore, within the time of the simulation. The binding at the intracellular site completely interrupted the water flux.

4.3 Discussion

In the current study, I have modeled the human isoform of the Aquaglyceroporin (AQP9). Repeated all-atom MD simulations indicated a structurally stable

model at the simulated timescale of 100 ns. By performing a series of simulations with the model of hAQP9 together with a set of recently identified specific and relatively potent inhibitors of mAQP9 [62], I found residues which were regularly contacted by the ligands. These residues were suggested for a mutational study. Point mutations of three (H82A, M91N and F180V) of the six suggested residues resulted in changes of binding affinity of already validated ligands. According to the homology model, these point mutation sites are located at the intracellular entrance of the water pore of hAQP9. Intracellular binding of an AQP inhibitor was also suggested by Migliati *et al.* [91].

Utilizing this information, I was able to identify a molecular docking approach for the identification of novel potent inhibitors. The application of this approach led to the identification of novel hAQP9 inhibitors from a chemical compound collection containing 1 million different structures. From 105 active compounds, 30 compound showed a significant inhibitory effect at 100 μM . From these, the top 18 (by activity) compounds were evaluated by titration response measurements. From these, the top 6 (by IC_{50}) compounds exhibit a half maximal inhibitory concentration of $\text{IC}_{50} = 3.6$ to 18 μM . These identified inhibitors (ID1–6) were structurally divergent from the originally identified substances [62] and therefore, provide additional starting points for future inhibitor optimizations. Efficacy studies with these inhibitors (ID1–6) on cells expressing the mutated hAQP9 isoforms provided additional evidence for the interaction of the identified inhibitors in close vicinity of the predicted interaction site. Computational support was provided by the significant reduction of the calculated p_f values of the hits. According to the poses of the hits at the target site, the hits penetrate the pore and span the NPA region, allowing contacts with Arg219 from the ar/R region, which is located at the most narrow region of the pore. Furthermore, free all-atom MD simulations with a validated inhibitor supported that the pore of hAQP9 is wide enough to allow the entering of drug-sized compounds. Notably, in 30 out of 40 simulations the partition of at least one compound into the lipid bilayer was observed. In order to bind at the intracellular site a ligand has to overcome the lipid bilayer. This suggests, that a sufficient lipophilicity is favorable for hAQP9 inhibitors. However, much longer time scales can be involved in drug binding than covered with this set of MD simulations as recently reported [17, 32, 114]. Therefore, effects that take place on the micro- and millisecond timescales may be considered in future studies.

A recently reported blind experimental screening against hAQP9 yield a hit rate of 3% [62]. In contrast the hit rate yield by the *in silico* supported screen reported here was 21%. Therefore, the hit rate was increased by 700% due to

the virtual screening. I note however, that the most potent substance identified with the described model was less effective than the best inhibitor identified in the blind small molecule screen (by 4–5 fold). It therefore remains to be clarified whether the *in silico* approach can be refined in order to identify more potent AQP9 inhibitors. So far, the optimization of the VS approach was only targeted to yield a large number of active compounds. In future studies, the proposed model of hAQP9 can be used to identify unsatisfied interaction sites between AQP9 and the know inhibitors or more favorable derivates of the inhibitors.

The pharmaceutical potential of AQP inhibitors has been emphasized in several review articles, most recently by Huber *et al.* [59], but relatively few organic AQP inhibitors have been described so far [15, 62, 88, 91, 95]. Specific and potent inhibitors are valuable for the study of AQP function, since they complement knockout and knockdown experiments, also when the compounds, as reported here, are not in perfect agreement with pharmacological requirements that are needed for an Investigational New Drug (IND) application. The general strategy outlined in this study is applicable to other members of the AQP family and related proteins.

4.4 Summary

Involvement of the AQP9 membrane channel has been suggested to be involved in the etiology of several diseases including *hyperglycemia* and type-2 diabetes [105]. By now, much of the knowledge about function of aquaporins has been explored using non-pharmacological methods as mRNA and protein expression studies and mutational studies. Also the association of mutations of genes encoding AQPs with human diseases has aided understanding of AQP function. However, a faster and more direct way to inhibit specific members of the AQP protein family would pave the way for a rigorous understanding of the function of aquaporins and their role in human pathologies. For that purpose, inhibitors that specifically inhibit certain members of the aquaporin family are invaluable tools for such investigations.

In the current study, novel potent inhibitors were identified. Furthermore, an atomic model of a the interaction site of human AQP9 for small compounds was established. The investigations presented herein, base on a composition of numerous experimental and computational methods as homology modeling, the combination of various molecular docking based strategies (blind docking, molecular docking with know compounds and virtual screening) as well as MD based methods (free all-atom MD simulations, p_f calculation). The putative

inhibitor binding site identified through this procedure was modified by site-directed mutagenesis and validated by fluorescence intensity assays. Using this model, it was possible to identify novel inhibitors of hAQP9 by a high-throughput virtual screening of 1 million compounds. The proposed model has been further supported by free MD simulations, in which the same ligand pose was generated that was suggested by molecular docking.

I greatly thank Dr. Michael Rützler who established the three independent cell lines for each mutated AQP9 isoform, performed the fluorescence assays and provided information about first inhibitors. In summary, we have successfully established a strategy for identification of small molecule inhibitors for hAQP9 that includes a 3D structural model of the putative binding site. In the future, the proposed molecules can be used in functional studies of AQP9 or serve as a basis for development of specific, sub-micromolar inhibitors. This computational approach can also be applied to other aquaporins and related proteins.

4.5 Methods

Molecular dynamics simulations were carried out using the Gromacs simulation software [53, 124]. The simulation box contained the protein tetramer embedded in a lipid bilayer of DMPC lipids and around 18,000 SCPE [7] water molecules. 78 potassium ions and 86 chloride ions were placed randomly in the simulation box corresponding to an ion concentration of $124 \mu\text{M}$. Ion parameters were taken from Dang [24]. Lipid parameters were taken from [9]. The whole system was simulated using the Amber ff99SB-ILDN [80] force field. Ligand parameters were calculated using the generalized amber force field (GAFF)[134]. ESP partial charges were calculated in a Hartree-Fock method at the 6-31G* level in agreement with the Amber ff99SB-ILDN force field. Restrained electrostatic potential (RESP) fitting was performed using Antechamber [133]. Long-range electrostatic interactions were calculated with the particle-mesh Ewald method [25, 36]. Short-range attractive van-der-Waals and repulsive Pauli interactions were described by a Lennard-Jones potential, which was cut off at 1 nm. The pressure was kept constant by coupling to a semi-isotropic Parrinello-Rahman barostat [98, 99] at 1 bar with a coupling constant of 5 ps. The Lincs algorithm [52] was used to constrain protein and lipid bond lengths and the Settle algorithm [92] was used to constrain all other bond lengths, allowing a time step of 2 fs. The temperature was kept constant by coupling the system to a velocity rescaling thermostat [8, 18] at 300 K with a coupling constant of 0.01 ps. The solvent and lipid molecules in simulation system were equilibrated for 10 ns before production. During this equilibration the atoms of the receptor were harmonically restrained with harmonic force constants of $1000 \text{ kJ}\cdot\text{mol}^{-1}\text{nm}^{-2}$. The simulation length of the production runs varied between 30 and 400 ns.

The calculation of the single channel osmotic water permeability coefficients p_f was done according to the collective diffusion model described by Zhu *et al.* [142]. The first 10 ns of the unrestrained simulations were removed for equilibration. The p_f values were computed individually for each subunit from the slope of the mean-square displacement (MSD) of the collective water coordinates. The slope of the MSD was approximated by the slope of a linear fitted line between 20 and 200 ps displacement time. According to 90 ns simulation time used in this analysis, the fit was carried out 450 times. For the analysis water molecules were used within a cylinder of length 1 nm and radius of 0.5 nm around the NPA motif. To exclude influence of the bulk, the z-position of the cylinder was adjusted to minimize the resulting p_f values. The reported values are mean values of the four computed p_f values from the individual subunits. The error bars correspond to

the respective standard deviation.

All structures targeted with molecular docking were energy minimized in the Amber ff99SB-ILDN [80] force field, applying 250 steps of the steepest decent algorithm implemented in Gromacs [53, 124]. Target structures for individual molecular docking algorithms used in this study were prepared as follows: The first molecular docking calculations were performed with Autodock-Vina [121] (here referred to as Vina; version 1.0.2). Input files of the target sites were generated using the Autodock-Vina plugin [113] for the PyMOL Molecular Graphics System (Version 1.3 from Schrödinger LCC). A standard grid spacing of 0.375 Å was used. The blind docking calculations were performed as described by Hetényi *et al.* [54, 55, 56] with Autodock-Vina (Vina), using whole hAQP9 tetramer as target site. Putative binding pockets were identified by visual inspection of the docking results. Conventional docking calculations were also performed with Vina. All atoms in a box of 21x21x28 Å³ around the intracellular pore entry were used as active site. An energy minimized homology model of hAQP9 and 39 representative structural clusters of the protein as well as 4 structures from equilibrated MD simulation containing reference ligands were used as target structures. Molecular docking with LeadIT [101] (version 2.0) (formerly FlexX) was done using the standard configuration. All atoms within a sphere of 10 Å around the intracellular pore entrance of hAQP9 served as target site. The high-throughput screen was performed using an energy minimized structure of the homology model of the hAQP9 wild type. A snapshot taken from an equilibrated MD simulation was utilized. The energy minimization was done in the Amber ff99SB-ILDN force field [80].

3D structures were prepared and hydrogenated with the program Conrina [107] (version 3.48). A structural screening database, containing 972307 drug-like compounds, was kindly provided by Enamine (<http://www.enamine.net>). Tanimoto-coefficients section were calculated using cheminformatics and machine learning software RDkit (<http://www.rdkit.org>) and default 2048 bit hash Daylight topological fingerprints (Section 2.1.4). The minimum path size was 1 bond, the maximum 7 bonds.

5

Identification of First Active Compounds

A question that has been raised regularly during the time of my studies was: “What should I do when no prior information of inhibitor molecules is available?”. In such cases it is not possible to benchmark and train virtual screening (VS) approaches beforehand as done in the former studies (Chapters 3 and 4). On the basis of the glutamine high-affinity transport system GlnPQ from *Lactococcus lactis* (*L. lactis*), the following study demonstrates a possible strategy for the identification of novel inhibitor molecules by molecular docking without prior optimization.

In the previous studies it was observed that a VS approach that performs well on one target does not necessarily perform well on other targets (compare FlexX and Vina in chapters 3 and 4). Therefore, I combined the prediction of putative novel inhibitors from different VS methods in order to increase the probability of identifying a successful screening approach. Accordingly, the primary goal of this study was the identification of a successful virtual screening approach rather than the identification of active compounds. This study resulted from a collaboration with Faizah Fulyani and Dr. Bert Poolman from the Groningen Biomolecular Science and Biotechnology Institute at the University of Groningen, who performed growth assays to test the suggested compounds for inhibitory activity.

5.1 Introduction

ATP-binding cassette (ABC) systems comprise a large super-family of membrane proteins which appear in archaea, eubacteria and eukarya [65]. The most prominent characteristic of these proteins is that they share a highly conserved ATPase domain, the ABC, which has been demonstrated to bind and hydrolyze adenosine triphosphate (ATP), thereby providing energy for a large number of biological processes [26]. There are three main classes of ABC-proteins: importers, exporters and a third class which does not involve transport but rather processes such as DNA repair and the translation or the regulation of gene expression.

ABC-transporters (importers and exporters) use ATP to translocate a broad number of substrates across the membrane including lipids, amino acids, peptides, proteins, metal ions, salt and hydrophobic compounds which also encompass drugs. A prominent member of the ABC-exporter family is the multidrug-resistance-protein (MDR1) that is known to cause resistance of tumors against a broad range of drugs [96].

The core structure of ABC-transporters consists of an assembly of two transmembrane (TM) subunits which build a pore and two intracellular ATP-binding subunits, together referred to as the *translocator*. Prokaryotic ABC transport systems involved in solute uptake (importers) employ an additional protein that captures the ligand and delivers it to the translocator. These so-called substrate-binding proteins (SBPs) are present in the periplasm of Gram-negative bacteria or they can be tethered to the membrane via a lipid or protein anchor (Gram-positive bacteria, *archaea*).

In 2002, it was discovered that members within the ABC class have the SBPs fused to the translocator domains [123], most notably the OTCN family involved in the uptake of osmoprotectants, taurine (alkyl sulfonates), cyanate and nitrate, and the PAO family which is specific for polar amino acids and opines (see ref. [26] for details on these families). In this case, it is referred to it as the substrate-binding domain (SBD). Within the PAO family, one or two domains are linked to the amino-terminal end of the TM subunit, and these are preceded by a signal sequence. Two of these chimeric substrate-binding/transmembrane proteins together with two ATP-binding cassettes form the functional unit for transport, and these systems thus have two or four substrate-binding domains [111].

ABC-transporters with SBDs attached to the transmembrane domain are present in (most prevalently) Gram-negative bacteria and Gram-positive bacteria including Gram-positive pathogens such as *Listeria monocytogenes*, *Staphylococcus au-*

reus, *Streptococcus pneumoniae*, *Enterococcus faecalis*, *Streptococcus pyogenes*. Many Gram-positive pathogens are multiple amino acid auxotrophs and require uptake of glutamine or glutamate via the Gln high affinity ABC-transport system GlnPQ that is predominantly expressed in Gram-positive bacteria [111]. Group B *streptococci*, for example, are the leading cause of neonatal sepsis and meningitis. In this pathogen, the GlnPQ transporter is implicated in the regulation of expression of fibronectin adhesins that are necessary for bacterial adhesion to other cells and therefore important for its virulence [118]. In *Lactococcus lactis* (*L. lactis*) GlnPQ is the only transporter for the uptake of glutamine (Gln) and glutamic acid (Glu) [111]. It has been shown that the disruption of the *glnPQ* gene leads to a loss of glutamine and glutamic acid uptake and cells do not longer grow on a medium with Gln or Glu as the sole source. However, the cells grow normally when Glu in the form of the dipeptide Ala-Gln was added and transported by a dipeptide transport system of *L. lactis*. Importantly, this class of ABC-transporters is not present in humans or other mammals, thus rendering GlnPQ a possible target for pharmaceutical treatment. Therefore, the interruption of Gln and Glu uptake by GlnPQ may be a possible route towards in the treatment against bacterial pathogens. However, no small molecule inhibitors of GlnPQ are currently known.

Recently, the structures of the substrate-binding domains (SBD1 and SBD2) from GlnPQ were resolved by X-ray crystallography [112]. Liganded structures were resolved to resolutions of 1.4 Å for SBD1 and 0.9 Å for SBD2. Furthermore, SBD2 was resolved in an unliganded open conformation at 1.5 Å resolution (Fig. 5.1) and one crystal structure was obtained from a tandem structure containing both SBD1 and SBD2 at 2.8 Å resolution. Although the two SBDs have a similar fold and similar binding pockets, the ligand affinity and specificity of SBD1 and SBD2 differ. For instance, the dissociation constant K_d (Section 2.1.1) for Gln is 91 μM on SBD1 and 0.9 μM on SBD2. At the primary structure level, SBD1 and SBD2 share 46% sequence identity. These structures enabled me to perform a structure based virtual screening (SBVS) for putative inhibitors of GlnPQ on the individual subunits. In order to identify a predictive virtual screening (VS) approach I combined a selection of compounds suggested from different VS programs Autodock-Vina [121] and FlexX [101] section targeting both SBD1 and SBD2. Afterwards, the predicted top compounds were purchased and evaluated in a functional assay bases on the growth rates of *L. lactis* in the presence of SBD1, SDBD2 or both. Adjacently, the most efficient VS approaches were deduced with the confirmed hits.

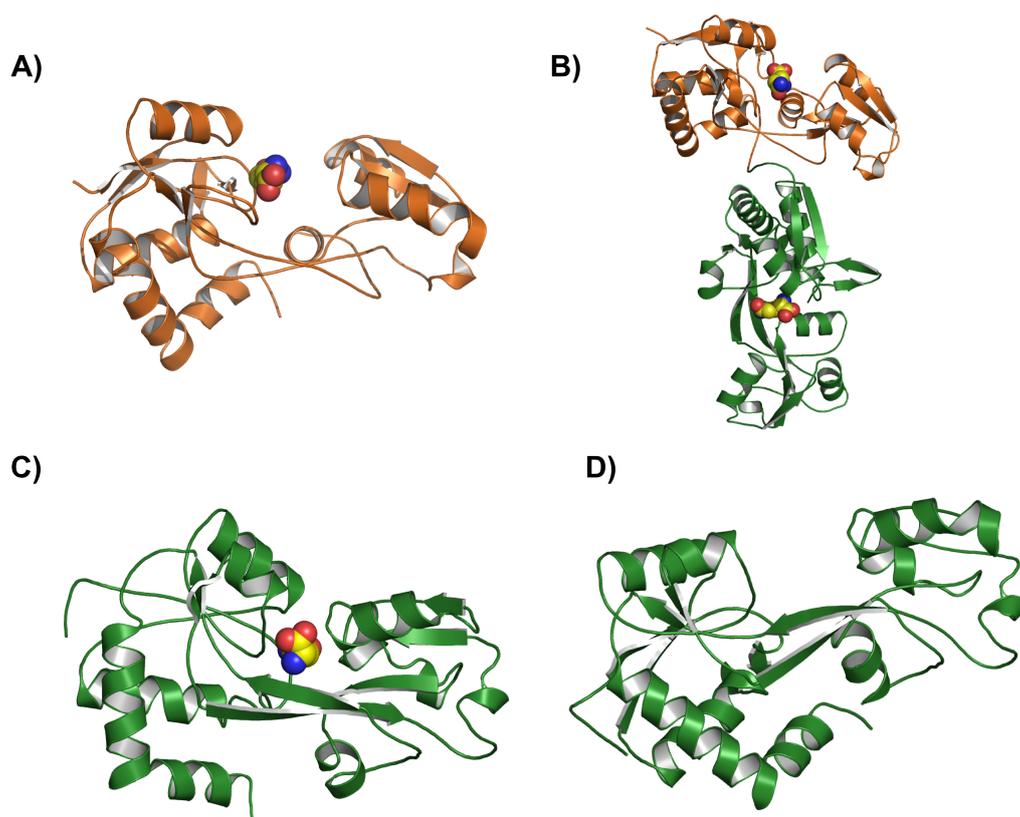


Figure 5.1: Obtained crystal structures. SBD1 liganded open (A), tandem SBD1/SBD2 open (B), SBD2 liganded closed (C) and unliganded open (D)

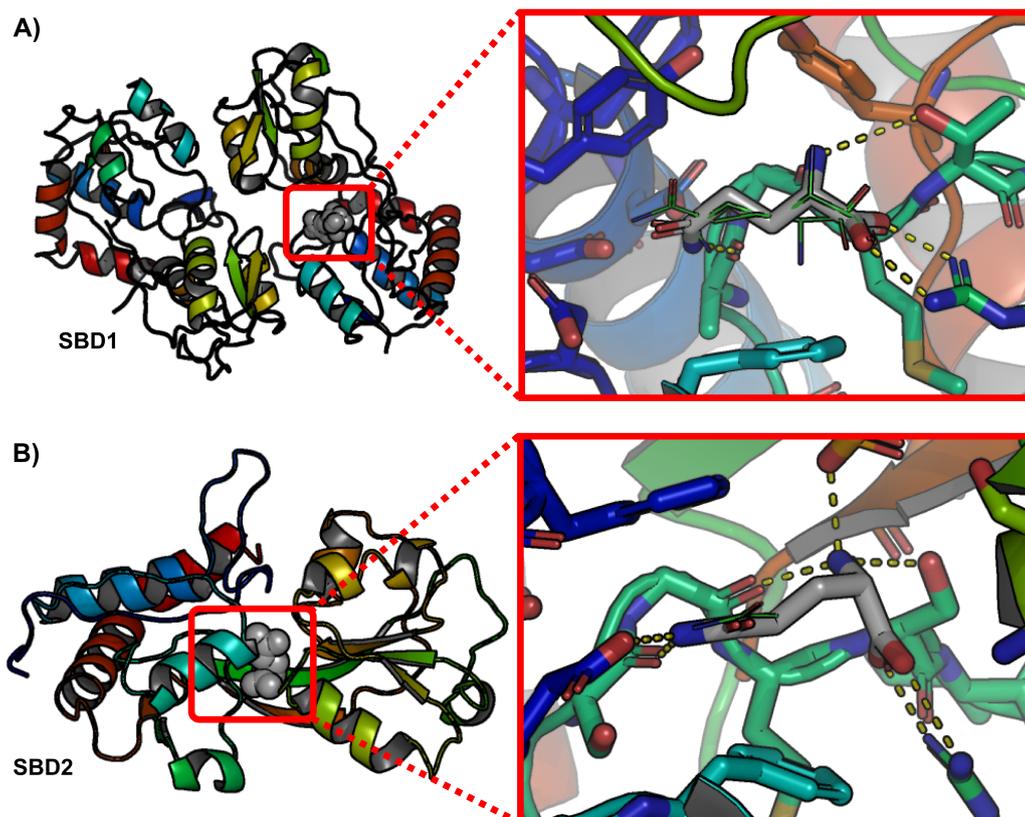


Figure 5.2: Crystal structures and docked poses (grey) of SBD1 (A) and SBD2 (B) using FlexX. The poses of glutamine as resolved in the crystal structure are shown in line representation. In the case of SBD2, the docking pose overlaps almost completely with the crystal structure.

5.2 Results

5.2.1 Reproduction of the Crystal Structure

No inhibitors of GlnPQ are currently available in the literature. Therefore, it was not possible to test and train a VS technique with respect to the enrichment of active compounds in advance. However, the crystal structures were available in the liganded forms (Fig. 5.1). This allowed a comparison of the docking poses of Gln with the conformations observed in the crystal structure as a first validation. The molecular docking program FlexX was used to dock the natural substrates of GlnPQ, namely glutamine and glutamic acid, against SBD1 and SBD2. The crystal structures of SBD1 and SBD2 with the highest resolutions were taken for the molecular docking calculation. All amino acids within a sphere of 6.5 Å radius around the substrate served as target sites. The molecular

Table 5.1: FlexX docking scores [$\text{kJ}\cdot\text{mol}^{-1}$] of glutamine, rank and RMSD [\AA] of the top poses to the substrate resolved in the crystal structure (by score and by RMSD) and corresponding experimentally obtained standard binding free energy ΔG_0^b [$\text{kJ}\cdot\text{mol}^{-1}$] derived from the K_d values [μM] at a temperature of 298K.

Criterion	Gln					Glu
	Score	Rank	RMSD	Exp. ΔG_0^b	Exp. K_d	Score
Lowest score SBD1	-29.7	1	2.4	-33	91	-18
Lowest score SBD2	-54.8	1	0.7	-45	0.9	-35
Lowest RMSD SBD1	-20.9	128	1.0	-33	91	n.e.
Lowest RMSD SBD2	-50.5	6	0.4	-45	0.9	n.e.

Table 5.2: Successfully screened compounds and intersections η , using the programs FlexX and Vina.

	SBD1	SBD2	η
FlexX	8,742	12,780	8,253
Vina	792,867	764,437	616,017
η	6,284	8,463	6,291

docking algorithm was able to generate the poses as found in the crystal structure (Figure 5.2). In the case of SBD2, the pose with the best overlap with the crystal structure was within the six top ranked poses. In the case of SBD1 the lowest docking scores for Gln were $-29.7 \text{ kJ}\cdot\text{mol}^{-1}$ for SBD1 and $-54.8 \text{ kJ}\cdot\text{mol}^{-1}$ for SBD2, and for Glu $-20.9 \text{ kJ}\cdot\text{mol}^{-1}$ for SBD1 and $-50.5 \text{ kJ}\cdot\text{mol}^{-1}$ for SBD2. For Gln, the tendency of these numbers agrees with the experiment ($-33 \text{ kJ}\cdot\text{mol}^{-1}$ for SBD1 and $-45 \text{ kJ}\cdot\text{mol}^{-1}$ for SBD2) [112](in preparation). However, the predicted scores deviate from the experimentally obtained values by 4–19 $\text{kJ}\cdot\text{mol}^{-1}$ (Tab. 5.1).

5.2.2 Structure Based Virtual Screening

In order to identify a VS approach for the identification of active compounds, I combined the prediction of four different methodically diverse VS approaches. Two different compound libraries were screened with FlexX and Vina. Library L_A contained 972,307 commercially available drug-like structures. Due to limitations of computational capacities a smaller but structurally more diverse Library

Table 5.3: Overview of the different strains of *L. lactis* involved in this study.

Strain-ID	Strain-Name	Available SBDs
S0	<i>L. lactis</i> Δ glnPQ	—
S12	<i>L. lactis</i> Δ glnPQ/pNZglnPQ	SBD1 and SBD2
S2	<i>L. lactis</i> Δ glnPQ/pNZglnPQ Δ SBD1	SBD2
S1	<i>L. lactis</i> Δ glnPQ/pNZglnPQ Δ SBD2	SBD1

L_B of drug-like compounds was used for FlexX. L_B contained 20,160 compounds. Library L_A was screened with Vina and Library L_B was screened with FlexX. The numbers of successfully screened compounds and the intersections are listed in Tab. 5.2. In order to screen for a proper VS approach, we applied 4 different methods (M1–M4) to select a set of around 100 compounds. These methods were:

- M1** 16 compounds were selected from the overlap of the best SBD1 or SBD2 compounds according the docking scores from FlexX and Vina
- M2** 20 compounds were selected by a consensus-score from the FlexX scores according to SBD1 **and** SBD2
- M3** 50 compounds were selected by a consensus-score of FlexX **and** Vina docking scores according to SBD1 **and** SBD2
- M4** 26 compounds were selected by best Vina scores according to SBD1 **or** SBD2

This selection resulted in a list of 104 compounds. These compounds (plus 2 predicted inactive compounds as negative controls) were purchased and evaluated in a functional assay to identify compounds that specifically inhibit GlnPQ as described in the following paragraph.

5.2.3 Experimental Validation

The GlnPQ transporter expressed in *L. lactis* was used to validate the *in silico* screening of transport inhibitors. *L. lactis* requires Glu (transported in the form of glutamic acid) or Gln for growth. A *L. lactis* Δ glnPQ strain (S0) was constructed that could be complemented *in trans* with wildtype glnPQ (S12) or used for testing of mutants in which either of the two SBDs were mutated or completely deleted. In addition, strains with cleaved SBD1 (S2) and cleaved SBD2

(S1) were constructed. For clarity, a summary of the reported strains herein is given in table 5.3.

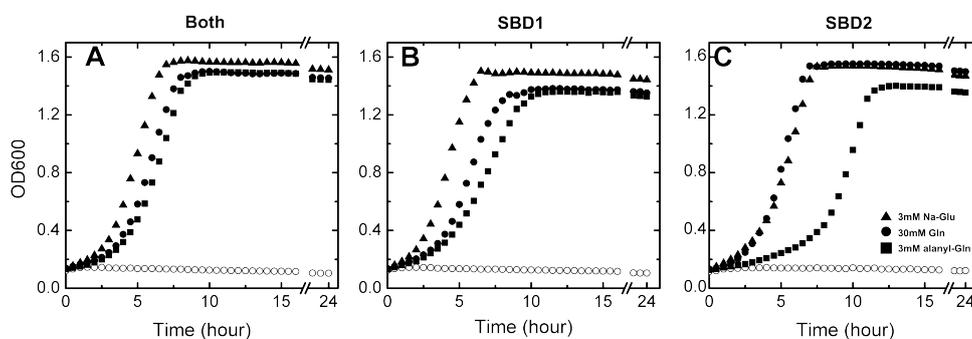


Figure 5.3: Growth of *L. lactis* strains S0 (A), S1 (B), S2 (C) in chemically defined medium (CDM) without Gln or Na-Glu (\circ), with 3 mM Na-Glu (\bullet), 30 mM Gln (\blacksquare), or 3 mM Ala-Glu (\blacktriangle). The cells were grown at 300 K in a 96 well plate format with a total cultivation volume of 300 μ l for at least 18 hours. The growth was followed by measuring the optical density at 600 nm at 30 minutes intervals.

Strain S12 (*L. lactis* Δ glnPQ/pNZglnPQ) did not grow in chemically-defined medium (CDM) lacking Gln or Glu (Fig. 5.3), neither strain S0 grew in CDM with Gln, showing that GlnPQ is essential for Gln transport and required for high rates of Gln uptake. Strain S12 grows at maximal rate provided 3 mM Na-Glu, 30 mM Gln or 3 mM alanyl-Gln present in the CDM (Fig. 5.3). Figure 5.4 shows that 0.1 mM of compound C1 does not inhibit growth of *L. lactis* strain S12 in the presence of Ala-Gln (panel A), whereas it does inhibit growth in the presence of Na-Glu (panel B) or Gln (panel C). At a concentration of 0.25 mM or higher, C2 is less specific and also affects growth in the presence of Ala-Gln. We note that growth inhibition can be manifested as a longer lag-phase, lower growth rate and/or lower cell yield (biomass formation), which may be due to adaptation of the cells to the stress imposed by the tested compounds. Although this complex behavior complicates the analysis of the inhibitors in terms of dose-response or IC_{50} values, the data allows rapid screening of compounds that specifically inhibit GlnPQ. The dipeptide alanyl-Gln (Ala-Gln) was taken up by the dipeptide and tripeptide transport protein (DtpT) [49]. The growth in the presence of Ala-Gln was used to discriminate between general inhibition by the tested compounds and GlnPQ-specific inhibition. The VS revealed at least 4 compounds (C1–C4)

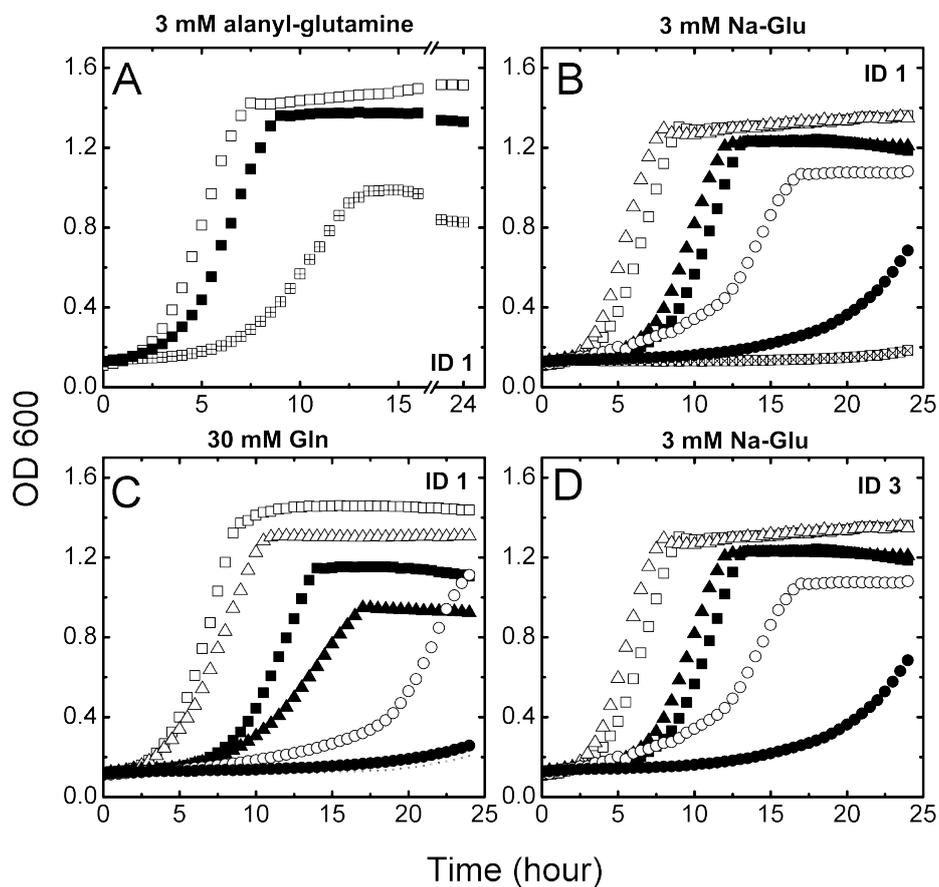


Figure 5.4: Panel A: Effect of compound C1 on the growth of *L. lactis* strain S12 in CDM in the presence of 3 mM Ala-Glu; no inhibitor (□), 0.1 mM compound C1 (■), and 0.25 mM compound C1 (⊠). Panel B, the same as panel A but with 3 mM Na-Glu instead of Ala-Glu. Panel C, the same as panel A but with 30 mM Gln instead of Ala-Glu. Symbols: strain S12 (□, ■), strain S1 (△, ▲), and strain S2 (○, ●). Panel D: Comparison of the effect of compound C1 on strain S12, S1 and S2. The inhibitors were dissolved in 100% (w/v) DMSO and diluted into the growth medium to the indicated concentrations and a final DMSO concentration of 2%.

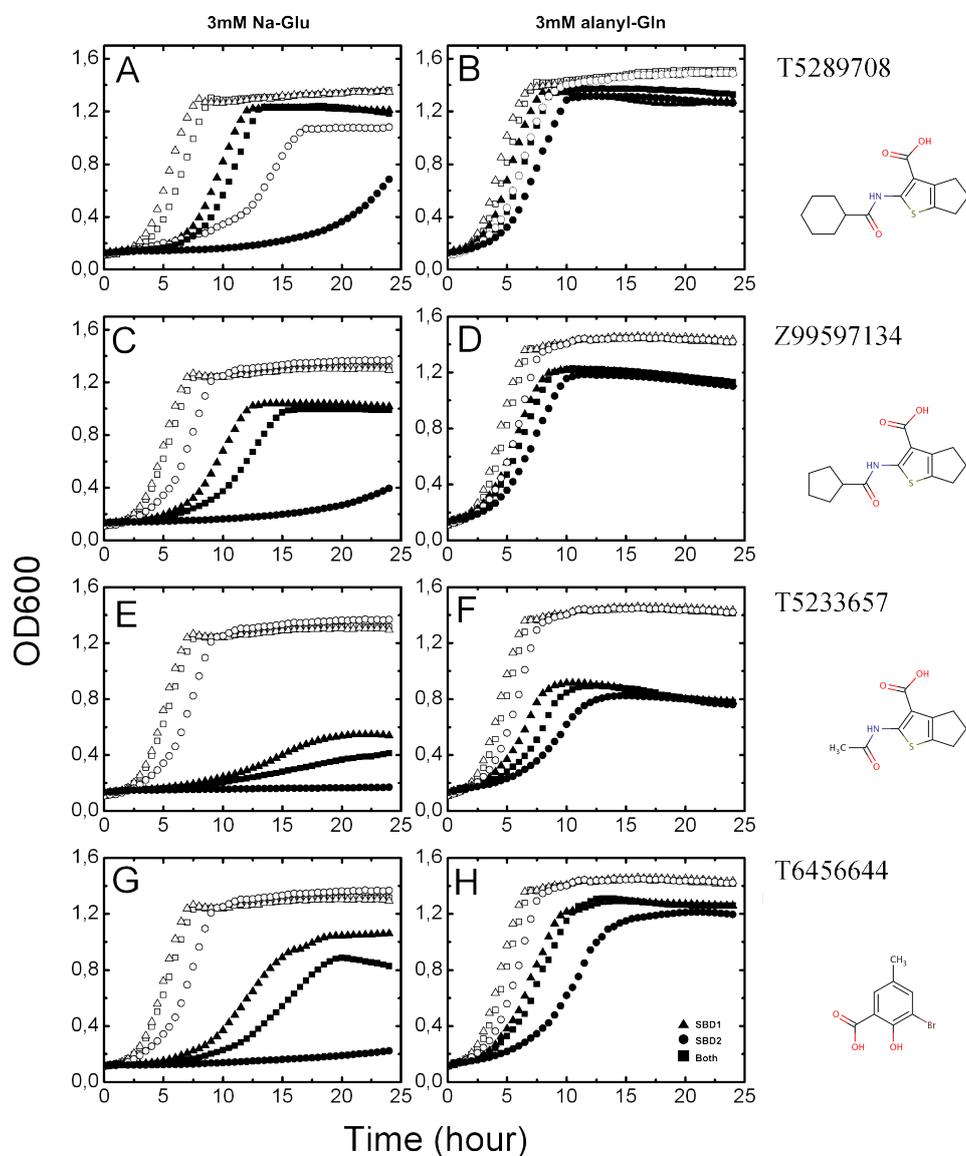


Figure 5.5: Comparison of related inhibitor compounds on the growth of *Lactococcus lactis* in the presence of 3 mM Na-Gln (left panel) and 3 mM Ala-Gln (right panel). Inhibition by C1 is shown in panels A and B at 0 and 0.1 mM for strains S12 (□, ■), S1 (△, ▲), and S2 (○, ●). The data for compound C2 at 0 and 0.25 mM are shown in panels C and D; data for compound C3 at 0 and 1 mM in panels E and F; and data for compound C4 at 0 and 1 mM in panels G and H. The chemical structures of the compounds are shown on the right of the panels.

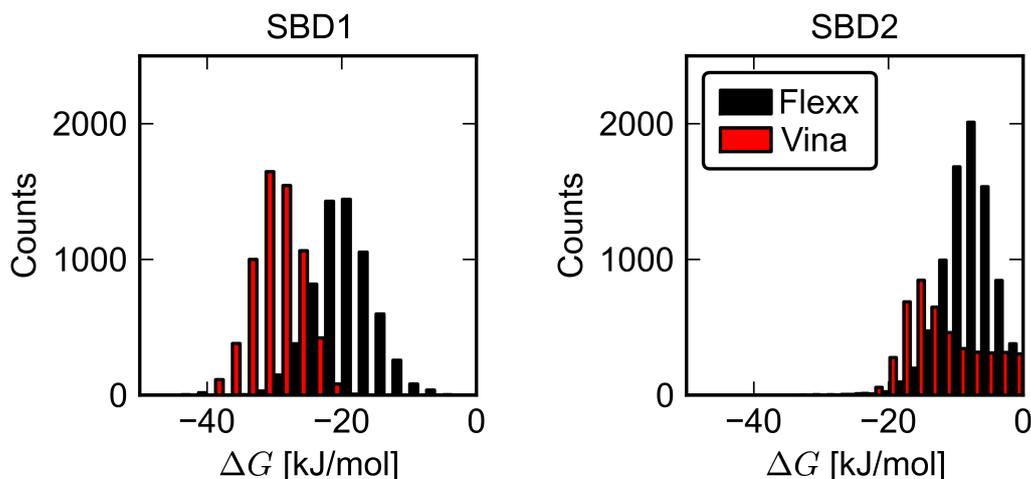


Figure 5.6: Distribution of the predicted binding free energies from the intersections in table 5.2 according to SBD1 (6,284 compounds) and SBD2 (8,463 compounds). Only compounds with scores below zero were considered.

that selectively inhibit GlnPQ. To determine whether compound C1 inhibits GlnPQ-mediated transport by interacting with SBD1 or SBD2, the inhibition of the wildtype GlnPQ (strain S12) was compared with that of strains S1 and S2. The effects of the compounds were exerted primarily on strain S2 which solely contained the SBD2 (Fig. 5.4, Panel D). Next, analogs of compound C1 were tested that have a five-ring (C2) or methyl group (C3) instead of six-ring structure linked to the core molecule. Figure 5.5 shows that the inhibition of growth and thus GlnPQ diminishes with decreasing size of this functionality. Specific inhibition was also observed for the unrelated structure C4, albeit at a relatively high concentration (Fig. 5.5).

5.2.4 Identification of First Inhibitors

The experimental assay confirmed at least four hits (C1–C4) that selectively inhibit GlnPQ. Three compounds (C1–C3) were identified using method M3 that is based on the diverse compound library L_B . Compounds C1–C3 share a common structure (2-methyl-4H,5H,6H-cyclopenta[b]thiophene-3-carboxylic acid). The comparison of the docking scores with the average values of all compounds which were evaluated with method M3 (including 6,291 compounds) (Tab. 5.4) and the distributions of the docking scores (Figure 5.6) reveals, that these compounds are scored lower than the average with FlexX on both SBD1 ($-19.3 \pm 4.3 \text{ kJ}\cdot\text{mol}^{-1}$) and SBD2 ($-7.7 \pm 3.9 \text{ kJ}\cdot\text{mol}^{-1}$), whereas the scores from Vina were moderate compared to the average value of $-29.9 \pm 3.6 \text{ kJ}\cdot\text{mol}^{-1}$ for SBD1. The average

Table 5.4: Predicted negative binding free energy ΔG by FlexX and Vina of the four hits according to SBD1 and SBD2 in units of $\text{kJ}\cdot\text{mol}^{-1}$. Average and standard deviations (STD) were calculated with respect to compounds from the intersections (Table 5.2).

ID	Enamine-ID	SBD1		SBD2		Method
		FlexX	Vina	FlexX	Vina	
C1	T5289708	-29.9	-30.5	-23.2	7.1	M3
C2	Z99597134	-29.2	-31.0	-23.0	-14.2	M3
C3	T5233657	-32.0	-26.8	-23.4	-10.5	M3
C4	T6456644	<i>n.e.</i>	-23.8	<i>n.e.</i>	-28.5	M4
Average		-19.3	-29.9	-7.7	<i>n.e.</i>	
STD		4.3	3.6	3.9	<i>n.e.</i>	

and standard deviation of the Vina scores according to SBD2 were not evaluated, since the distribution deviated from a normal distribution. However, the distribution of the Vina scores according to SBD2 shows a peak at $-17\text{ kJ}\cdot\text{mol}^{-1}$. On the basis of an approximated standard deviation of $4\text{ kJ}\cdot\text{mol}^{-1}$, as observed for the other distributions, C4 is scored significantly lower than this reference value at SBD2. For compounds C1–C3, both FlexX and Vina predict higher binding free energies at SBD1 compared to SBD2. Notably, only for C4, a higher predicted binding free energy at SBD2 was predicted. Considering the poses generated by Vina according to SBD2, Figure 5.7 illustrates that compounds C3 and C4 were placed inside the cavity showing significant overlap with the reference pose of the substrate. Compounds C1 and C2 were placed at the surface of SBD2 similar to the poses that were generated by FlexX (Fig. 5.8). The FlexX poses of C1–C3 are very similar. The aromatic ring systems of C1–C3 are placed in a small hydrophobic pocket. Furthermore, all three molecules show interaction with an Arg residue. Both Vina and FlexX on average suggest higher affinities with respect to SBD1 than to SBD2. However, a comparison of the scores that were obtained by the hits with the average values shows, that all four hits C1–C4 were selected on the basis of the scores of SBD2.

5.3 Discussion

The glutamine high-affinity ABC-transport system GlnPQ is predominantly expressed in Gram-positive bacteria [111] and provide the main route for fast import of Gln and Glu. Several gram-positive pathogens are amino acid auxotrophs

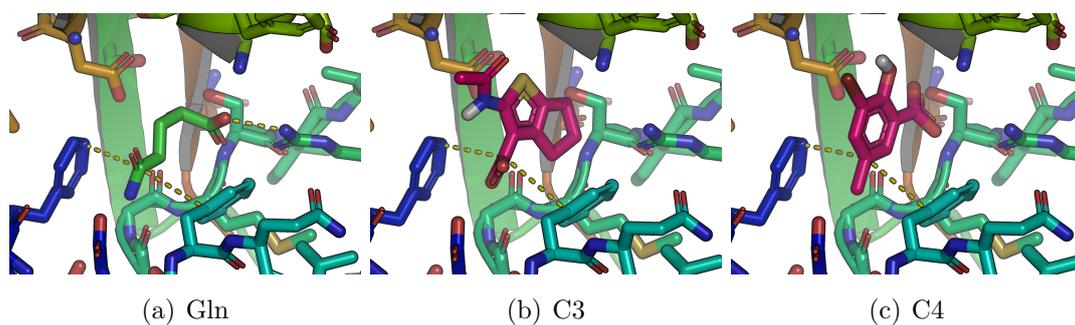


Figure 5.7: Reference pose of glutamine in the crystal structure and the docking poses of compounds C3 and C4 generated by Vina which overlapped with the Gln resolved in the crystal structure of SBD2.

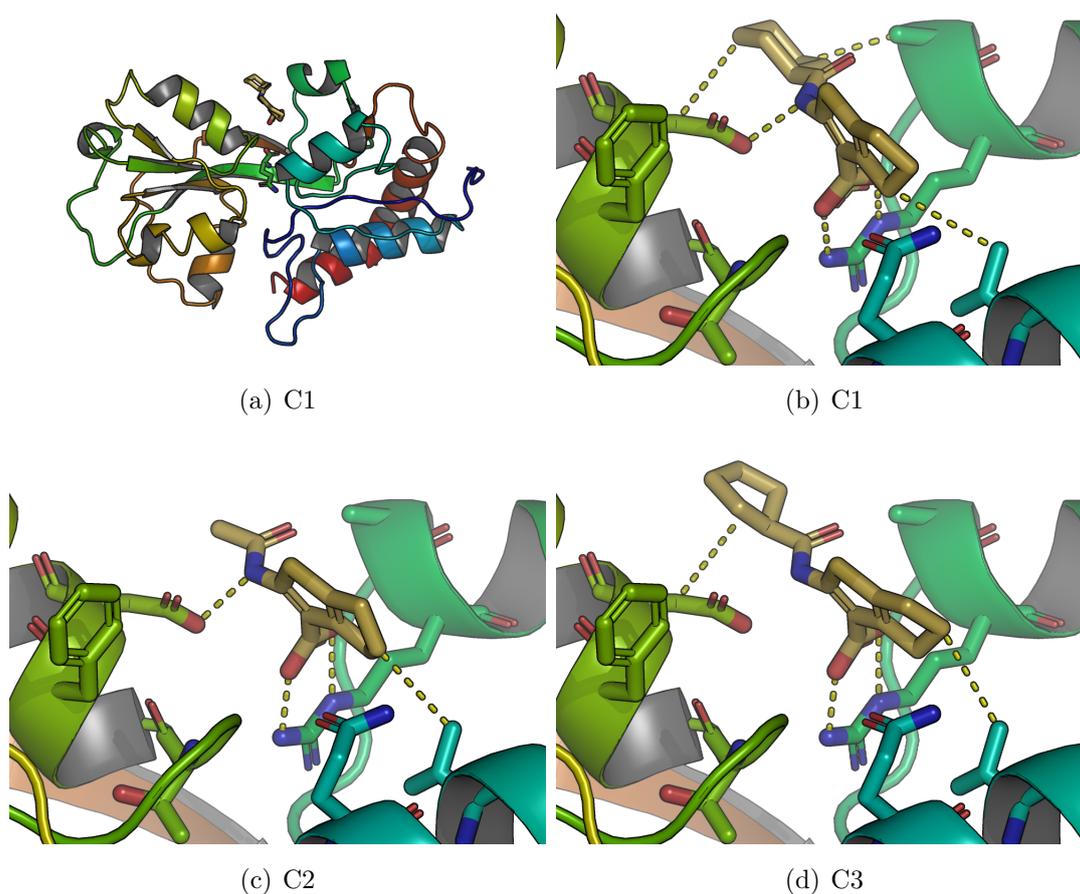


Figure 5.8: Reference pose of glutamine in the crystal structure (A) and the docking poses of C1(B), C3(C) and C2(D) generated with FlexX.

and require uptake of Gln or Glu via GlnPQ for growth. Therefore, they need GlnPQ for growth. Importantly, this class of ABC-transporters with two substrate binding domains (SBD1 and SBD2) is not present in humans or other mammals, rendering GlnPQ a possible target for pharmaceutical treatment.

This study presents the identification of the first selective inhibitors of GlnPQ by subsequent virtual screening (VS) and functional assaying. The combined prediction according to four methodically diverse VS approaches (M1–M4) led to the identification of at least 4 selective inhibitors of GlnPQ. Growth experiments with four different strains (S0,S12,S1,S2) of *L. lactis* indicate that the effects of the most thoroughly analyzed compounds (C1–C4) originate from interactions with SBD2. The inhibition could be due to competition with glutamate or glutamine binding or the transfer of these ligands from the SBD2 to the TM domain.

The sampling of the crystallographic pose of the substrate Gln with FlexX supported a successful screening approach. The predicted standard binding free energy of $-54.8 \text{ kJ}\cdot\text{mol}^{-1}$ for Gln at SBD2 is in the range of the experimental value of $-45 \text{ kJ}\cdot\text{mol}^{-1}$ obtained from the dissociation constant $K_d=0.9 \mu\text{M}$ at 298 K. At SBD1 the docking score at SBD1 ($-29.7 \text{ kJ}\cdot\text{mol}^{-1}$) is closer to the experimentally obtained value of $-33 \text{ kJ}\cdot\text{mol}^{-1}$ ($K_d=91 \mu\text{M}$ at 298 K). The predicted values for Glu are at least $10 \text{ kJ}\cdot\text{mol}^{-1}$ smaller compared to Gln. Generally, molecular docking cannot be expected to yield a quantitative agreement to experimentally determined affinities due to the severe underlying assumptions.

The identification of first hits was an important achievement that makes it possible to analyse the predictive power of the individual VS approaches. Within a set of 106 compounds at least 4 specifically inhibit GlnPQ. Compounds C1–C3 were predicted by method M3 that is a consensus of four docking scores according to SBD1 and SBD2 and the programs FlexX and Vina. Therefore, the method is designed to predict compounds which are predicted to interact with both SBD1 and SBD2 with both programs. However, the experimental results show that the compounds perform their inhibitory effect by interacting with SBD2. This is also reflected in the docking scores of the individual programs, since the scores on SBD2 for C1–C3 are significantly lower than the average of the FlexX scores on SBD1. The Vina scores of compounds C1–C3 are rather moderate, on both SBD1 and SBD2 and therefore the FlexX scores accounted for the selection of these compounds. However, method M2 did not identify active compounds even though this method consists of a consensus of FlexX scores according to SBD1 and SBD2. An analysis of the consensus scores of C1–C3 according to method M2 revealed that these compounds appeared in the first 50 ranks of M2 (data

not shown). Therefore, the extension to the top 50 compounds in M2 would have led to an identification of these compounds, too.

Compound C4 was selected by method M4, that is purely based on Vina scores. Here, no consensus method was applied, since compounds were selected which were top scored either on SBD1 or SBD2 or on both sides. The docking scores of C4 show that this compound was selected due to the low score on SBD2.

On SBD2, a significant amount of compounds was too large to fit into the small binding pocket and therefore Vina assigned positive scores to these compounds. However, on average both programs predicted a higher affinity to SBD1 compared to SBD2, what was not directly in agreement with the experimental results. This may be due to the selected conformation of SBD2. With around $35 \text{ kJ}\cdot\text{mol}^{-1}$ the highest predicted binding free energies were around $20 \text{ kJ}\cdot\text{mol}^{-1}$ above the value for the natural substrate Gln. It is important to stress that we used the crystal structure of the liganded closed conformation of SBD2 for the docking. Therefore, all side chains were in an optimal conformation for Gln, but not for other compounds explaining the low score of Gln. Further analysis should clarify whether conformational changes in SBD2 increase the average predicted binding free energy. The next step could be to target the crystal structure of the liganded or unliganded open conformation of SBD2 in a VS.

The low number of active compounds makes an analysis with respect to the most powerful screening technique difficult. To clarify this, a post-screening with the hits and some decoys could be used to identify the method with the best enrichment. The coupling of the individual scoring functions by the consensus in methods M1-M4 complicates the deduction of the predictive factors. For future studies it may be more useful to select compounds based on individual molecular docking calculations and the subsequent construction of a proper consensus approach. Alternatively, method M1 could be modified to take the union of the best compounds instead of than the intersection. In general, it seems that both FlexX and Vina have predictive power on SBD2.

After this test, a high-throughput virtual screening (HTVS) using both FlexX and Vina will be applied in order to identify more inhibitors. In parallel, compounds C1-C4 will be evaluated on GlnPQ homologues from other (pathogenic) species and the inhibition characterized in terms of IC_{50} in order to estimate their binding affinity.

5.4 Methods

The target sites of GlnPQ were generated with the graphical user interface of FlexX (LeadIT version 2.0) using the liganded crystal structures of SBD1 and SBD2 with the highest resolution. All residues within a sphere of 6.5 Å around the substrate were defined as target site. Standard parameters were used for weights of the scoring function and the number of intermediate solutions for each fragment. Input files were generated using the AutoDock plug-in [113] for the program PyMOL [28] using the liganded crystal structures of SBD1 and SBD2 with the lowest resolution. A cubic box of 7x7x7 Å³ centered around the substrate defined the target sites.

Two structural screening databases of commercially available compounds, were kindly provided by Enamine (<http://www.enamine.net>). Library L_A contained a broad set of 972,307 drug-like compounds, whereas library L_B contained a diverse set of 20,160 drug-like compounds. 3D structures were prepared and protonated with the program Conrina [106] (version 3.48). Tanimoto-coefficients section were calculated using cheminformatics and machine learning software RDkit (<http://www.rdkit.org>) and default 2048 bit hash Daylight topological fingerprints (Section 2.1.4). The minimum path size was 1 bond, the maximum 7 bonds.

Experimental Setup

Growth experiments were performed in the 96-well format, using a total cultivation volume of 300 µl. Exponentially growing pre-cultures were used to inoculate the 300 µl chemically defined medium (CDM), supplemented with a 10,000-fold dilution of *L. lactis* NZ9700 spent medium (containing the inducer, nisin A), 5 µg/ml chloramphenicol, and 1% (w/v) glucose [10]. In total 106 putative inhibitors (purchased from Enamine Ltd, Ukraine) were screened; the inhibitors were dissolved in 100% (w/v) DMSO and diluted into the growth medium to a final DMSO concentration of 2%. At 2% DMSO, growth of *L. lactis* 9000, carrying plasmids with wildtype or mutant derivatives of glnPQ, is not yet affected. Cell growth was followed for at least 18 hours at 30 minutes interval by measuring the optical density at 600 nm, using an automated microtiterplate reader (Biotek).

6

Outlook: Hit-Optimization based on Molecular Docking

Once a list of predicted compounds has been compiled with supporting data and hit compounds are confirmed, the selection process to prioritize chemical series for the hit-to-lead follow up begins. A commonly used technique at this stage is *hit evolution* where hit derivatives are generated in order to find more active and selective compounds [70]. This chapter reports on the generation of a molecular docking based algorithm for the structure based optimization of active compounds and discusses the first results obtained during the development process. The findings presented herein provide an outlook for future studies.

6.1 Introduction

Compared to the estimated number of drug-like compounds in the drug-like chemical space (10^{60})[13], the scope of compounds in a typical drug-like virtual screening library is vanishingly small (10^5 to 10^7). Therefore, the probability to find compounds with more favorable physiochemical properties is higher when the chemical space outside the boundaries of the library is taken into account. A virtual screening (VS) approach that was trained on a particular target and that already has successfully identified novel active compounds provides evidence of

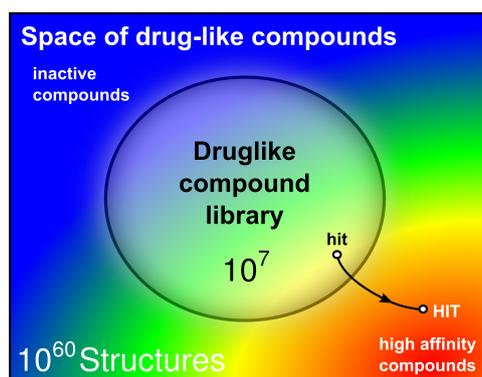


Figure 6.1: A typical drug-like compound library used for virtual screening is significantly small compared to the druglike chemical space. Using known active structures in combination with a validated virtual screening approach it possibly even more favorable compounds can be identified.

predictive power concerning the experimental activities. Hence, the development of an optimization algorithm for the generation of promising derivatives based on FlexX was performed. Therefore, such a VS approach may also be promising with respect to the directed optimization of active compounds towards higher activity and selectivity. Chapter (4) provided a successful VS study where active compounds in the low micromolar range were identified. The compound ID1 with the an IC_{50} of $= 3.1 \mu\text{M}$ (Fig. 4.10) was subjected to an iterative optimization procedure aiming to identify more efficient and more active derivatives. According to the results in chapter 4, the optimization was performed based on the FlexX total score.

6.2 Implementation and Results

6.2.1 Compound Modification

The first of the three essential modules for the optimization is the modification of the input structure or the structures, if multiple input structures are provided. The modification algorithm can either be used to freely explore the surrounding in the chemical space of a parent structure or be implemented in a way that only child structures are generated that fall into certain constraints. The former option would allow the exploration of undesired but potentially fruitful chemical space, since e.g. the inconsistency with the rule of five does not necessarily mean undesired biochemical properties. The latter option focuses on a chemical sub-space with a higher probability for success and therefore may save computational

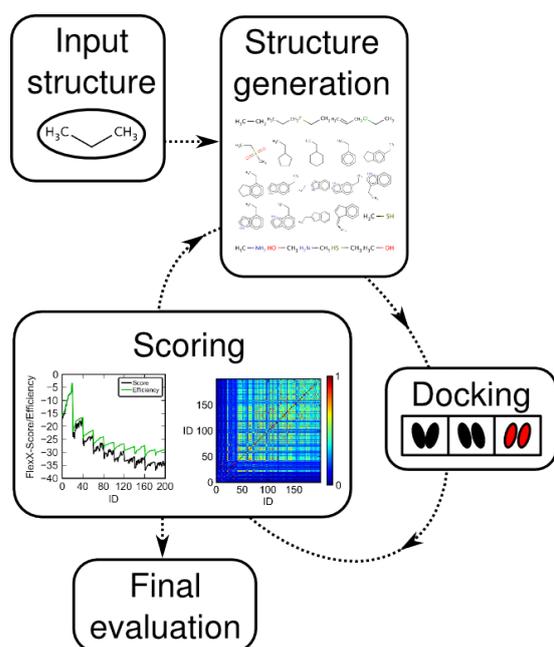


Figure 6.2: Scheme of the optimization algorithm: input structures are the initial seed for the modification module. The generated derivatives are docked to one or more receptor structures. The efficiency is evaluated based on the docking results and structure of the derivatives. The top N scored compounds serve as seed for the next iteration until the efficiency is converged. Finally, all top scored compounds are reevaluated by a second (more accurate) scoring procedure.

time and experimental effort. The generated set of 3D molecular structures is then processed to the next module to estimation their affinities.

6.2.2 The Docking Module

The second module reflects the affinity estimation of the modified structures on the target receptor. Here it is referred to as “docking” since the molecular docking algorithm that was employed to predict the hits identified in chapter 4 is applied. The docking part is the same as described in section ?? since it has been validated against a set of known active compounds and it successfully predicted novel active compounds. Multiple receptor structures can be targeted either as target or anti-target. After docking the molecule to each structure a consensus-score is calculated. Here, the average of the target structures is taken as final docking score.

Initially, the optimization was performed on the basis of the pure docking score. As a first test, compound ID1 (Fig. 6.3) was subjected to the optimization algorithm. The top 20 compounds of each iteration served as seed for the next iteration. As illustrated in Fig. 6.3 the docking score decreased during the iterations to a value of approximately $-50 \text{ kJ}\cdot\text{mol}^{-1}$. Also structural convergence was observed as indicated by the similarity of the generated structures (Fig. 6.3, Panel C). The optimized structures of this first generation optimization algorithm were enriched with hydrogen acceptors and donors forming favorable interactions with the receptor and achieving low scores (data not shown). In the following,

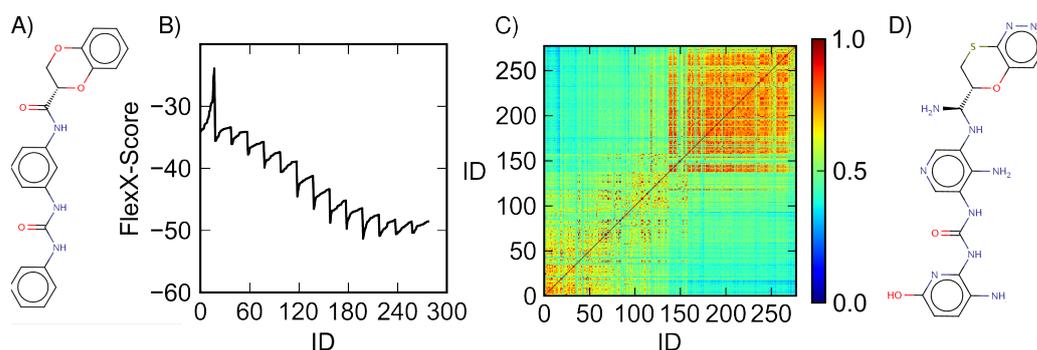


Figure 6.3: Optimization of compound ID1 (A) (Chapter 3). The FlexX total score (B) used for optimization converged to a value of $-50 \text{ kJ}\cdot\text{mol}^{-1}$. The similarity of the generated structures indicates structural convergence (C). The optimized structures were endowed with several hydrogen donors (D).

the optimization according to an efficiency function that takes the structure of the ligand into account is discussed.

6.2.3 Final Scoring

A final evaluation was inserted to increase the number of true-positives. All compounds which were used as seeds during the optimization procedure were reevaluated. Technically, this was done by rescoring poses with a second scoring function, namely the Hyde module that is implemented in the LeadIT/FlexX software suite [101] and which is designed for reduction of false positives. Hereby, the poses which were generated with FlexX were optimized in the Hyde force field. The relaxed pose is then scored with the Hyde scoring function which is supposed to describe hydrogen bonding and desolvation effects more accurately than the FlexX scoring function [110].

6.2.4 Compound Efficiency

The pure docking score optimizes for high receptor-ligand interactions and therefore may leave the drug-like chemical space. Therefore, an efficiency E was introduced, that based on both the docking score and properties of the ligand. The efficiency penalized deviations from certain reference values in order to guide the process to both active and preferential drug-like properties. The most efficient structures then serve as seed in the next iteration. A first approach to ensure drug-like structure was to penalize deviations from a target compound mass of $m_0 = 400 \text{ Da}$ and a target fraction carbon atoms $f_0 = 57\%$. These values were

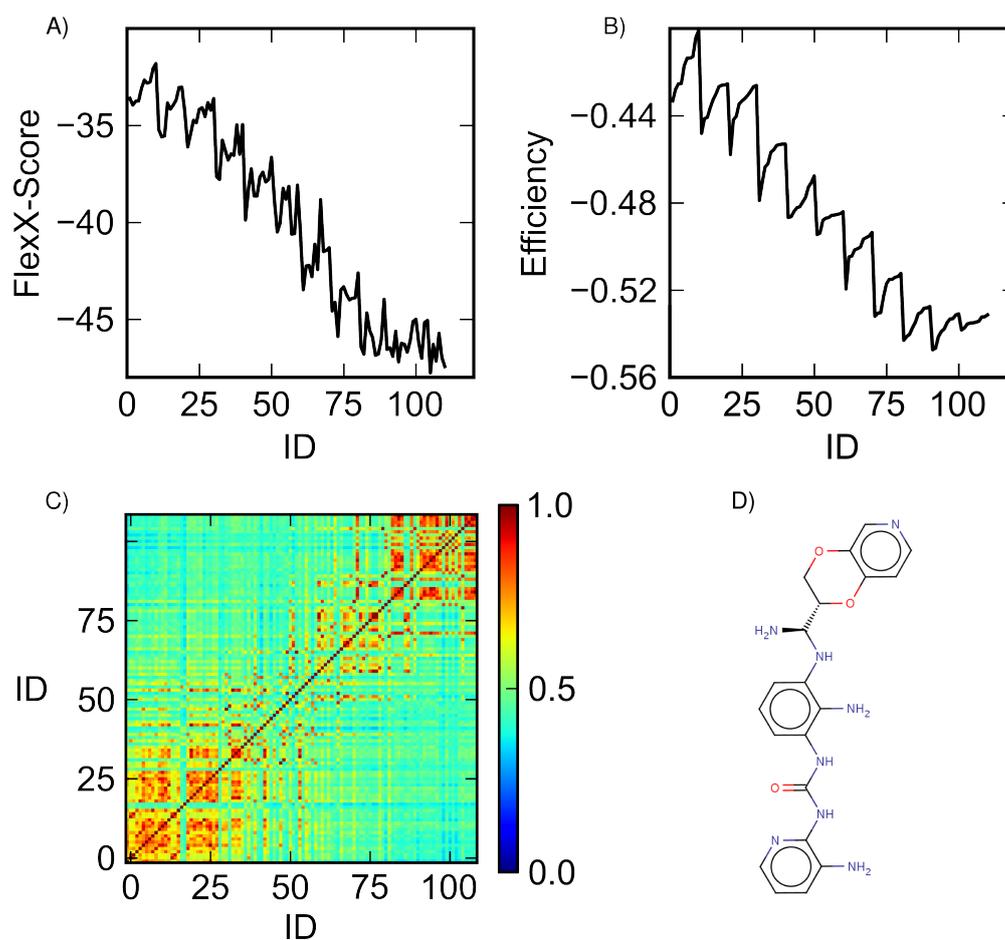


Figure 6.4: Results obtained by using the efficiency E for the optimization of ligand ID1 from chapter 4. Both docking score (A) and efficiency (B) converge very fast. The similarity of the top scored compounds (C) indicates structural convergence. The final structure reevaluated with Hyde (D).

taken to calculate the ligand efficiency:

$$E = S \cdot e^{-\frac{(m-m_0)^2}{\sigma_m}} \cdot e^{-\frac{(f-f_0)^2}{\sigma_f}} \quad (6.1)$$

where S is the FlexX total score, m the mass, f the fraction of carbon atoms and the respective target values m_0 and f_0 . For σ_m and σ_f define the strength of the restraining and were set to 300 Da and 30 % respectively to allow deviations from the target value. The restrained mass restricts the size of the ligands. Without mass restrictions the size of the ligands is only limited by the size of the target site. The restricted fraction of carbon atoms implicitly penalizes inefficient interactions. The optimization run was repeated using E for the optimization. In this case only the top 10 compounds of each docking run served as seed for the next iteration. The algorithm converged in only 11 iterations to a FlexX score of $-47 \text{ kJ}\cdot\text{mol}^{-1}$ (Fig. 6.4). The optimized structure had less hydrogen bond acceptors and donors than the optimized structure using the pure FlexX total score (Fig. 6.4). Notably, both compounds were similar in size. The current implementation of the optimization algorithm only allows additions and replacements of atoms but basically no deletions. Therefore, the compounds can only grow upon a particular core structure. In order to achieve even more efficient structures it would be necessary to allow deletions.

6.2.5 Quasi-*de novo* Design

The current implementation only allows growth of the compounds upon a certain core structure. This limits the chemical space that can be covered by the optimization algorithm. One possibility to overcome this restriction could be to allow atomic deletions and therefore shrinking and regrowth of the compounds. Alternatively to deletions, it is possible to use smaller fragments of the initial compound. A first attempt to explore this possibility was done using the urea motif that was observed to be present in various hAQP9 inhibitors (Fig. 4.10) and also present in compound ID1. Using the urea motif as input, the docking score converged to a value of $-50 \text{ kJ}\cdot\text{mol}^{-1}$ (Fig. 6.5). The resulting top scored structure was structurally different compared to the previous top scored compounds. This attempt revealed the possibility to use this approach for the generation of completely new scaffolds by using very small initial structures, therefore providing a smooth transition to the *de novo* design of novel inhibitors.

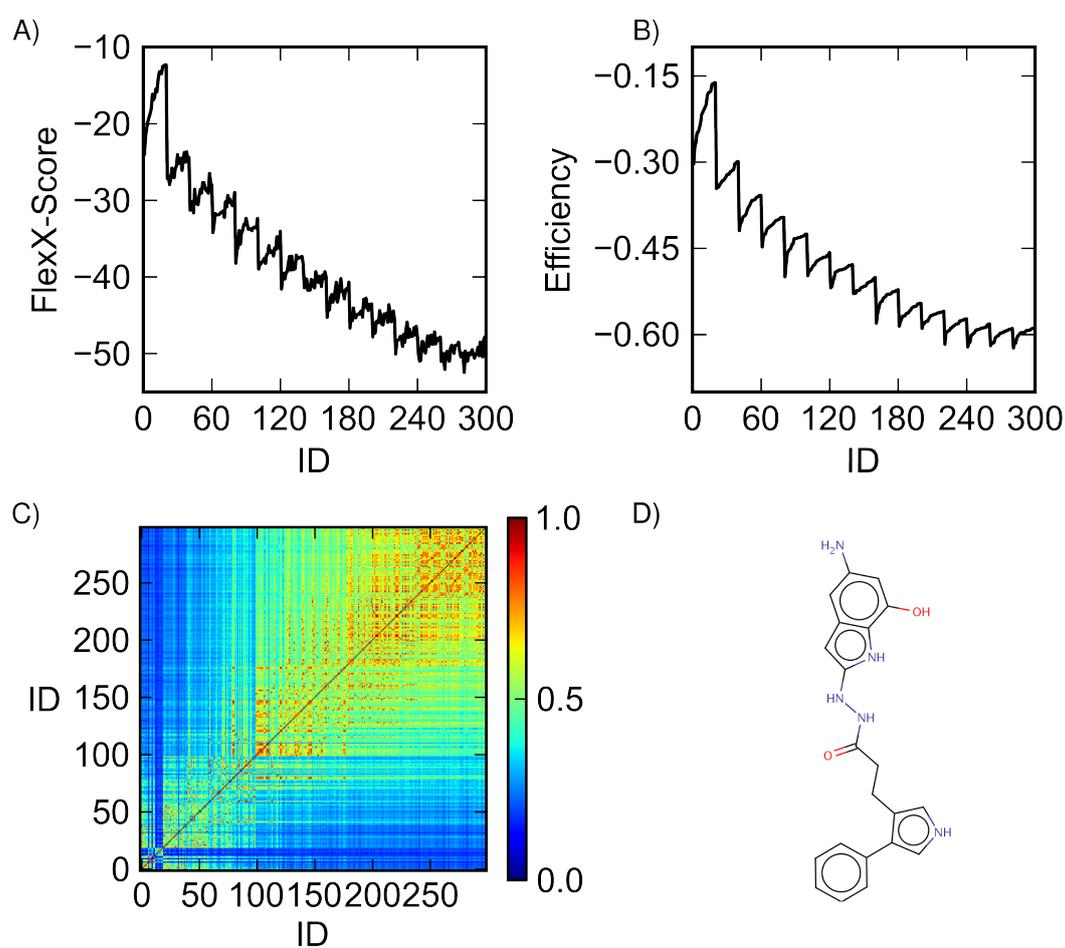


Figure 6.5: Results obtained from the urea motif. Both the docking score (A) and the efficiency (B) reach comparable level as in the previous cases (6.3 and 6.4). The similarity of the top scored compounds (C) indicates structural convergence. The final structure reevaluated with Hyde (D).

6.3 Outlook

First outcomes for the future optimization of active compounds were obtained using a molecular docking based iterative algorithm: The optimization with respect to the FlexX total score leads to molecules with non-drug-like properties (e.g. mass, number of hydrogen bond donors and acceptors). The weighting of the total score by exponential weighting functions can partially compensate for this deficiency. Larger molecules are restricted to a chemical subspace which is defined by the core structure of the compound and the modification module. Smaller fragments which serve as initial seed structure generate more diverse structures which are unrelated to the original active compound, except for the common scaffold. The use of very small fragments or even single atoms for the initial seed may be useful for the *de novo* design of compounds. However, the predictive power of these approach remains to be tested. Physiological experiments have to clarify the extend of dimension of structural changes that is most effective for the optimization towards or the design of novel high-affinity ligands. Furthermore, the generation of modified structures needs to be biased towards the synthesizable chemical space. At the current state, this algorithm is in an experimental and preliminary state and is not technically matured. In order keep the optimization algorithm well-arranged and flexible it was concepted and implemented in modules which can independently be developed or exchanged by more sophisticated algorithms in future studies. For example, the scoring can be replaced by other methods with (preferentially) enhanced predictive power as e.g. consensus scoring. The same holds for the final evaluation which could be replaced by relative or absolute binding free energy calculations based on simulation techniques.

6.4 Methods

Tanimoto-coefficients section were calculated using the cheminformatics and machine learning software RDkit (<http://www.rdkit.org>) and default 2048 bit hash Daylight topological fingerprints (Section 2.1.4). The minimum path size was 1 bond, the maximum 7 bonds. The target value for the fraction of carbon atoms was obtained from the screening library from the Maybridge catalogue of drug-like compounds (<http://www.maybridge.com>) containing 56,213 structures. The fraction of carbon atoms of the compounds in this library was $57.3 \pm 10.36\%$. Molecular docking with LeadIT [101] (version 2.0) (formerly FlexX) was done using a non standard configuration. The maximum number of solutions per iteration and the maximum number of solutions per fragmentations were set to 2000. All atoms within a sphere of 10 Å around the intracellular pore entrance of hAQP9 served as target site for the molecular docking. The high-throughput screen was performed using an energy minimized structure of the homology model of the hAQP9 wild type. A snapshot taken from an equilibrated MD simulation was utilized. The energy minimization was done in the Amber ff99SB-ILDN force field [80].

The structural modification was implemented in five steps:

- Conversion to SMILES strings
- Modification of SMILES strings
- Filter with respect to solubility
- Sorting and removal of redundant and unphysical/unstable structures
- Generation of 3D structures

The conversion to SMILES strings was done with the program Open Babel (<http://openbabel.org>). The generated strings were modified by adding and replacing groups. For that purpose a set of chemical groups was defined that was added to carbon atoms “C” and “c” (aromatic carbons) or the end of a chemical rest groups. These groups were in particular:

C, =C, =O, O, N, F, Cl, =F,

c1cccc1,

NC(=O),

C1=CC=CC=C1,

C1C=CC=C2CCCC12,

C1C=CC2=C1C=CC=C2,

C1=CC2=C(C1)C=CC=C2,
C1=CC2=C(C=CC2)C=C1,
C1=CC2=C(CC=C2)C=C1,
C1=CC=CC2=C1CC=C2,
C1=CC=CC2=C1C=CC2,
c1ccc2cc[nH]c2c1,
c1c[nH]c2ccccc12,
c1ccc2CCc2c1,
c1ccc2[nH]ccc2c1,
c1cccc2[nH]ccc12,
c1cccc2cc[nH]c12,
C1CCCC1,
C1CCCCC1,
=C1C=Cc2ccccc12,
C1=CC2=C(C1)C=CC=C2, C1=CC2=C(N1)C=CC=C2,
c1cc2ccccc2[nH]1, c1cc2ccccc2c1

Furthermore, a set of sites for atomic exchanges was defined (Tab. 6.1).

Table 6.1: Atoms as defined in the “exchange site” column were replaced by alternative atoms (including bond types). For each exchange site a set of alternative atoms defined.

Exchange site	Alternative groups
C	O, N, S
F	N, O, C, Cl
CC	C(C)
O	C, N, S, Cl
N	O, C
c	c(C), n, o, s
Cl	F, O, N, C
o	n, c
n	o, c(N), c
s	c, o
=O	H, O, N, =N, =C, =S

7

Conclusions

In this thesis I explored the potential of molecular docking and molecular dynamics simulations for the development of small molecule inhibitors of membrane channel proteins. The four studies presented herein cover the identification of a reliable virtual screening approach solely on the basis of a crystal structure (Chapter 5), the optimization of a virtual screening approach (Chapter 3) and the modeling of the binding process (Chapter 4). Finally the optimization of inhibitors by molecular docking (Chapter 6) is addressed. Hereby, the choice of the computational methods used in these studies was oriented towards the experiments which were performed in parallel by collaborators. The two major findings are that conventional molecular docking can gain considerable hit rates when optimized for membrane channel proteins and that the binding of inhibitors of hAQP9 takes place at the intracellular site.

The optimization of molecular docking was addressed in chapter 3. Based on an optimized consensus approach that combined the predictions of four molecular docking programs, I established 14 novel inhibitors of $K_V1.1-(1.2)_3$ with affinities down to the sub-micromolar range. These compounds were found to inhibit the current carried by $K_V1.1-(1.2)_3$ channels by more than 80 % at 10 μ M. Compared to blind experimental screenings this is an improvement in enrichment of two to three orders of magnitude. Furthermore, two of these compounds exhibited at least 30-fold higher potency in inhibition of $K_V1.1-(1.2)_3$ over a set of cardiac

ion channels (hERG, Nav1.5, and Cav1.2), indicating a pronounced selectivity for $K_V1.1-(1.2)_3$ and therefore meet a first set of cardiac safety constants. It is important to note that the final selection of compounds was purely based on computationally obtained values. Therefore, these results represent the pure predictive power of the algorithm. The consensus approach developed herein can be applied to other targets in particular other potassium channels.

The inhibition of the human water channel protein Aquaporin 9 (hAQP9) was addressed in chapter 4. On the basis of a sequence homology to the glycerol facilitator (GlpF) a structural model of hAQP9 was established. Simulations with known murine AQP9 inhibitors suggested residues putatively involved in ligand binding. Residues at the intracellular site of hAQP9 were confirmed to be involved in the binding process by mutagenesis and subsequent fluorescence assays performed by collaborators. The intracellular site was then targeted by a molecular docking based virtual screening for the identification of putative inhibitors. The rate of active compounds was increased by a factor 7 compared to a small blind screening. The activities of the most thoroughly analysed compounds lie in the low micromolar range. When docked to the intracellular site the calculated p_f values were significantly decreased compared to simulations without ligands. Importantly, both the simulations of spontaneous binding of a known active compound in a free all-atom MD simulation and experimental data supported intracellular binding and confirm the binding pose identified by molecular docking.

Comparing the studies in chapters 3 and 4, the results show that the predictive power of an individual molecular docking program strongly depends on the target structure, since the performance of both Vina and FlexX was opposite on both targets $K_V1.1-(1.2)_3$ and hAQP9. Furthermore, the results obtained in the $K_V1.1-(1.2)_3$ study indicate that the virtual screening approach used in the hAQP9 study can be optimized further with respect to the identification of more active inhibitors.

Chapters 5 and 6 report on ongoing studies and cover the initial phase for the identification of a successful virtual screening approach and the hit optimization phase when active compounds are known. The identification of the first active compounds of GlnPQ allows to optimize a virtual screening approach for this particular target in the future. The active compounds were selected on the basis of substrate binding domain SBD2 in agreement with the experiment. However, on average higher affinities were predicted on SBD1 indicating the need for structural refinement of the target structure SBD2.

The approach presented in chapter 6 was designed in order to optimize known

inhibitors based on validated scores from molecular docking with validated predictive power. The use of the pure docking score of FlexX however led to non-drug-like molecules. Therefore, the balancing of the scoring with respect to the docking scores and molecular properties (solubility, weight, etc.) or structural features (e.g. the fraction of carbon atoms) in order to obtain drug-like structures is currently studied. So far, no predictions of optimized compounds have been tested experimentally. This approach may also be applied for the *de novo* design of active compounds.

Bibliography

- [1] P. Agre, L. S. King, M. Yasui, W. B. Guggino, O. P. Ottersen, Y. Fujiyoshi, A. Engel, and S. Nielsen. Aquaporin water channels—from atomic structure to clinical medicine. *J Physiol*, 542(Pt 1):3–16, Jul 2002.
- [2] M. And er, V. B. Luzhkov, and J. Aqvist. Ligand binding to the voltage-gated Kv1.5 potassium channel in the open state—docking and computer simulations of a homology model. *Biophys J*, 94(3):820–831, Feb 2008.
- [3] E. J. Arroyo, Y. T. Xu, L. Zhou, A. Messing, E. Peles, S. Y. Chiu, and S. S. Scherer. Myelinating schwann cells determine the internodal localization of Kv1.1, Kv1.2, Kvbeta2, and Caspr. *J Neurocytol*, 28(4-5):333–347, 1999.
- [4] J. C. Baber, W. A. Shirley, Y. Gao, and M. Feher. The use of consensus scoring in ligand-based virtual screening. *J Chem Inf Model*, 46(1):277–288, 2006.
- [5] J. Baldwin, C. H. Michnoff, N. A. Malmquist, J. White, M. G. Roth, P. K. Rathod, and M. A. Phillips. High-throughput screening for potent and selective inhibitors of plasmodium falciparum dihydroorotate dehydrogenase. *J Biol Chem*, 280(23):21847–21853, Jun 2005.
- [6] E. Beitz, B. Wu, L. M. Holm, J. E. Schultz, and T. Zeuthen. Point mutations in the aromatic/arginine region in aquaporin 1 allow passage of urea, glycerol, ammonia, and protons. *Proc Natl Acad Sci U S A*, 103(2):269–74, 2006.
- [7] H. J. C. Berendsen, J. Grigera, and T. Straatsma. The missing term in effective pair potentials. *J Phys Chem*, 91(24):6269–6271, 1987.
- [8] H. J. C. Berendsen, J. P. M. Postma, W. F. van Gunsteren, A. DiNola, and J. R. Haak. Molecular dynamics with coupling to an external bath. *J Chem Phys*, 81:3684–3690, 1984.

- [9] O. Berger, O. Edholm, and F. Jähnig. Molecular dynamics simulations of a fluid bilayer of dipalmitoylphosphatidylcholine at full hydration, constant pressure, and constant temperature. *Biophys J*, 72(5):2002–2013, May 1997.
- [10] R. P.-A. Berntsson, N. A. Oktaviani, F. Fusetti, A.-M. W. H. Thunnissen, B. Poolman, and D.-J. Slotboom. Selenomethionine incorporation in proteins expressed in *Lactococcus lactis*. *Protein Sci*, 18(5):1121–1127, May 2009.
- [11] R. B. Best and G. Hummer. Optimized molecular dynamics force fields applied to the helix-coil transition of polypeptides. *J Phys Chem B*, 113(26):9004–9015, Jul 2009.
- [12] K. H. Bleicher, H.-J. Böhm, K. Müller, and A. I. Alanine. Hit and lead generation: beyond high-throughput screening. *Nat Rev Drug Discov*, 2(5):369–378, May 2003.
- [13] R. S. Bohacek, C. McMartin, and W. C. Guida. The art and practice of structure-based drug design: a molecular modeling perspective. *Med Res Rev*, 16(1):3–50, Jan 1996.
- [14] H. J. Böhm. The development of a simple empirical scoring function to estimate the binding constant for a protein-ligand complex of known three-dimensional structure. *J Comput Aided Mol Des*, 8(3):243–256, Jun 1994.
- [15] H. L. Brooks, J. W. Regan, and A. J. Yool. Inhibition of aquaporin-1 water permeability by tetraethylammonium: involvement of the loop E pore region. *Mol Pharmacol*, 57(5):1021–1026, May 2000.
- [16] S. Bruckner and S. Boresch. Efficiency of alchemical free energy simulations. II. improvements for thermodynamic integration. *J Comput Chem*, 32(7):1320–1333, May 2011.
- [17] I. Buch, T. Giorgino, and G. D. Fabritiis. Complete reconstruction of an enzyme-inhibitor binding process by molecular dynamics simulations. *Proc Natl Acad Sci U S A*, 108(25):10184–10189, Jun 2011.
- [18] G. Bussi, D. Donadio, and M. Parrinello. Canonical sampling through velocity rescaling. *The Journal of Chemical Physics*, 126:014101, 2007.

- [19] C. A. Chang, W. Chen, and M. K. Gilson. Ligand configurational entropy and protein binding. *Proc Natl Acad Sci U S A*, 104(5):1534–1539, Jan 2007.
- [20] P. S. Charifson, J. J. Corkery, M. A. Murcko, and W. P. Walters. Consensus scoring: A method for obtaining improved hit rates from docking databases of three-dimensional structures into proteins. *J Med Chem*, 42(25):5100–5109, Dec 1999.
- [21] S. K. Coleman, J. Newcombe, J. Pryke, and J. O. Dolly. Subunit composition of Kv1 channels in human CNS. *J Neurochem*, 73(2):849–858, Aug 1999.
- [22] M. Congreve, C. W. Murray, and T. L. Blundell. Structural biology and drug discovery. *Drug Discov Today*, 10(13):895–907, Jul 2005.
- [23] G. Crooks. Nonequilibrium measurements of free energy differences for microscopically reversible markovian systems. *Journal of Statistical Physics*, 90:1481–1487, 1998. 10.1023/A:1023208217925.
- [24] L. X. Dang. Mechanism and thermodynamics of ion selectivity in aqueous solutions of 18-crown-6 ether: A molecular dynamics study. *J Am Chem Soc*, 117:6954–6960, 1995.
- [25] T. Darden, D. York, and L. Pedersen. Particle mesh Ewald: An $N \log(N)$ method for Ewald sums in large systems. *J Chem Phys*, 98:10089–10092, 1993.
- [26] E. Dassa and P. Bouige. The ABC of ABCS: a phylogenetic and functional classification of ABC systems in living organisms. *Res Microbiol*, 152(3-4):211–229, 2001.
- [27] N. Decher, B. Pirard, F. Bundis, S. Peukert, K.-H. Baringhaus, A. E. Busch, K. Steinmeyer, and M. C. Sanguinetti. Molecular basis for Kv1.5 channel block: conservation of drug binding sites among voltage-gated K⁺ channels. *J Biol Chem*, 279(1):394–400, Jan 2004.
- [28] W. DeLano. The PyMOL Molecular Graphics System. *DeLano Scientific San Carlos, CA, USA*, 2002.
- [29] J. Devaux, M. Gola, G. Jacquet, and M. Crest. Effects of K⁺ channel blockers on developing rat myelinated CNS axons: identification of four types of K⁺ channels. *J Neurophysiol*, 87(3):1376–1385, Mar 2002.

- [30] R. Dias and W. F. de Azevedo. Molecular docking algorithms. *Curr Drug Targets*, 9(12):1040–1047, Dec 2008.
- [31] D. A. Doyle, J. M. Cabral, R. A. Pfuetzner, A. Kuo, J. M. Gulbis, S. L. Cohen, B. T. Chait, and R. MacKinnon. The structure of the potassium channel: molecular basis of K⁺ conduction and selectivity. *Science*, 280(5360):69–77, Apr 1998.
- [32] R. O. Dror, A. C. Pan, D. H. Arlow, D. W. Borhani, P. Maragakis, Y. Shan, H. Xu, and D. E. Shaw. Pathway and mechanism of drug binding to g-protein-coupled receptors. *Proc Natl Acad Sci U S A*, 108(32):13118–13123, Aug 2011.
- [33] M. D. Eldridge, C. W. Murray, T. R. Auton, G. V. Paolini, and R. P. Mee. Empirical scoring functions: I. The development of a fast empirical scoring function to estimate the binding affinity of ligands in receptor complexes. *J Comput Aided Mol Des*, 11(5):425–445, Sep 1997.
- [34] J. Eldstrom and D. Fedida. Modeling of high-affinity binding of the novel atrial anti-arrhythmic agent, vernakalant, to Kv1.5 channels. *J Mol Graph Model*, 28(3):226–235, Oct 2009.
- [35] P. Ertl. Cheminformatics analysis of organic substituents: identification of the most common substituents, calculation of substituent properties, and automatic identification of drug-like bioisosteric groups. *J Chem Inf Comput Sci*, 43(2):374–380, 2003.
- [36] U. Essmann, L. Perera, M. L. Berkowitz, T. Darden, H. Lee, and L. G. Pedersen. A smooth particle mesh Ewald method. *J Chem Phys*, 103:8577–8592, 1995.
- [37] N. Eswar, D. Eramian, B. Webb, M.-Y. Shen, and A. Sali. Protein structure modeling with MODELLER. *Methods Mol Biol*, 426:145–159, 2008.
- [38] T. Fawcett. An introduction to ROC analysis. *Pattern Recognit Lett*, 27(8):861–874, 2006. z.
- [39] M. Feher. Consensus scoring for protein-ligand interactions. *Drug Discov Today*, 11(9-10):421–428, May 2006.
- [40] R. A. Fenton, H. B. Moeller, S. Nielsen, B. L. de Groot, and M. Rutzler. A plate reader-based method for cell water permeability measurement. *Am J Physiol Renal Physiol*, 298(1):F224–30, 2010.

- [41] R. A. Friesner, J. L. Banks, R. B. Murphy, T. A. Halgren, J. J. Klicic, D. T. Mainz, M. P. Repasky, E. H. Knoll, M. Shelley, J. K. Perry, D. E. Shaw, P. Francis, and P. S. Shenkin. Glide: a new approach for rapid, accurate docking and scoring. 1. method and assessment of docking accuracy. *J Med Chem*, 47(7):1739–1749, Mar 2004.
- [42] D. Fu, A. Libson, L. J. Miercke, C. Weitzman, P. Nollert, J. Krucinski, and R. M. Stroud. Structure of a glycerol-conducting channel and the basis for its selectivity. *Science*, 290(5491):481–486, Oct 2000.
- [43] M. K. Gilson, J. A. Given, B. L. Bush, and J. A. McCammon. The statistical-thermodynamic basis for computation of binding affinities: a critical review. *Biophys J*, 72(3):1047–1069, Mar 1997.
- [44] M. K. Gilson and H.-X. Zhou. Calculation of protein-ligand binding affinities. *Annu Rev Biophys Biomol Struct*, 36:21–42, 2007.
- [45] T. Gonen, P. Sliz, J. Kistler, Y. Cheng, and T. Walz. Aquaporin-0 membrane junctions reveal the structure of a closed water pore. *Nature*, 429(6988):193–7, 2004.
- [46] A. D. Goodman, T. R. Brown, J. A. Cohen, L. B. Krupp, R. Schapiro, S. R. Schwid, R. Cohen, L. N. Marinucci, A. R. Blight, and F. M.-F. S. Group. Dose comparison trial of sustained-release fampridine in multiple sclerosis. *Neurology*, 71(15):1134–1141, Oct 2008.
- [47] A. D. Goodman, T. R. Brown, K. R. Edwards, L. B. Krupp, R. T. Schapiro, R. Cohen, L. N. Marinucci, A. R. Blight, and M. S. F. Investigators. A phase 3 trial of extended release oral dalfampridine in multiple sclerosis. *Ann Neurol*, 68(4):494–502, Oct 2010.
- [48] A. D. Goodman, T. R. Brown, L. B. Krupp, R. T. Schapiro, S. R. Schwid, R. Cohen, L. N. Marinucci, A. R. Blight, and F. M.-F. Investigators. Sustained-release oral fampridine in multiple sclerosis: a randomised, double-blind, controlled trial. *Lancet*, 373(9665):732–738, Feb 2009.
- [49] A. Hagting, E. R. Kunji, K. J. Leenhouts, B. Poolman, and W. N. Konings. The di- and tripeptide transport protein of *Lactococcus lactis*. A new type of bacterial peptide transporter. *J Biol Chem*, 269(15):11391–11399, Apr 1994.

- [50] T. A. Halgren, R. B. Murphy, R. A. Friesner, H. S. Beard, L. L. Frye, W. T. Pollard, and J. L. Banks. Glide: a new approach for rapid, accurate docking and scoring. 2. Enrichment factors in database screening. *J Med Chem*, 47(7):1750–1759, Mar 2004.
- [51] L. Heginbotham, Z. Lu, T. Abramson, and R. MacKinnon. Mutations in the K⁺ channel signature sequence. *Biophys J*, 66(4):1061–1067, Apr 1994.
- [52] B. Hess, H. Bekker, H. J. C. Berendsen, and J. G. E. M. Fraaije. LINCS: A linear constraint solver for molecular simulations. *J Comput Chem*, 18:1463–1472, 1997.
- [53] B. Hess, C. Kutzner, D. Van Der Spoel, and E. Lindahl. GROMACS 4: Algorithms for highly efficient, load-balanced, and scalable molecular simulation. *Journal of chemical theory and computation*, 4(3):435–447, 2008.
- [54] C. Hetényi and D. van der Spoel. Efficient docking of peptides to proteins without prior knowledge of the binding site. *Protein Science*, 11(7):1729–1737, 2002.
- [55] C. Hetényi and D. van der Spoel. Blind docking of drug-sized compounds to proteins with up to a thousand residues. *FEBS Lett*, 580(5):1447–1450, Feb 2006.
- [56] C. Hetényi and D. van der Spoel. Toward prediction of functional protein pockets using blind docking and pocket search algorithms. *Protein Sci*, 20(5):880–893, May 2011.
- [57] V. Hornak, R. Abel, A. Okur, B. Strockbine, A. Roitberg, and C. Simmerling. Comparison of multiple Amber force fields and development of improved protein backbone parameters. *Proteins*, 65(3):712–725, Nov 2006.
- [58] S.-Y. Huang, S. Z. Grinter, and X. Zou. Scoring functions and their evaluation methods for protein-ligand docking: recent advances and future directions. *Phys Chem Chem Phys*, 12(40):12899–12908, Oct 2010.
- [59] V. J. Huber, M. Tsujita, and T. Nakada. Aquaporins in drug discovery and pharmacotherapy. *Mol Aspects Med*, online:–, 2012.
- [60] T. J-F and B. C. I. Evaluating virtual screening methods: good and bad metrics for the "early recognition" problem. *J Chem Inf Model*, 47(2):488–508, 2007.

- [61] C. Jarzynski. Nonequilibrium equality for free energy differences. *Phys Rev Lett*, 78:2690–2693, 1997.
- [62] S. Jelen, S. Wacker, C. Aponte-Santamaria, M. Skott, A. Rojek, U. Johanson, P. Kjellbom, S. Nielsen, B. L. de Groot, and M. Rutzler. Aquaporin-9 protein is the primary route of hepatocyte glycerol uptake for glycerol gluconeogenesis in mice. *J Biol Chem*, 286(52):44319–25, 2011.
- [63] M. O. Jensen, D. W. Borhani, K. Lindorff-Larsen, P. Maragakis, V. Jogini, M. P. Eastwood, R. O. Dror, and D. E. Shaw. Principles of conduction and hydrophobic gating in K⁺ channels. *Proc Natl Acad Sci U S A*, 107(13):5833–5838, Mar 2010.
- [64] Y. Jiang, A. Lee, J. Chen, M. Cadene, B. T. Chait, and R. MacKinnon. Crystal structure and mechanism of a calcium-gated potassium channel. *Nature*, 417(6888):515–522, May 2002.
- [65] P. Jones and A. George. The ABC transporter structure and mechanism: perspectives on recent research. *Cellular and Molecular Life Sciences*, 61(6):682–699, 2004.
- [66] W. L. Jorgensen. The many roles of computation in drug discovery. *Science*, 303(5665):1813–1818, Mar 2004.
- [67] W. L. Jorgensen, J. Chandrasekhar, J. D. Madura, R. W. Impey, and M. L. Klein. Comparison of simple potential functions for simulating liquid water. *J Chem Phys* 79, 79:926, 1983.
- [68] W. L. Jorgensen, D. Maxwell, and J. Tirado-Rives. Development and testing of the OPLS all-atom force field on conformational energetics and properties of organic liquids. *J Am Chem Soc*, 118(45):11225–11236, 1996.
- [69] S. I. V. Judge and C. T. Bever. Potassium channel blockers in multiple sclerosis: neuronal Kv channels and effects of symptomatic treatment. *Pharmacol Ther*, 111(1):224–259, Jul 2006.
- [70] G. M. Keseru and G. M. Makara. Hit discovery and hit-to-lead approaches. *Drug Discov Today*, 11(15-16):741–748, Aug 2006.
- [71] J. G. Kirkwood. Statistical mechanics of fluid mixtures. *J Chem Phys*, 3:300, 1935.

- [72] D. B. Kitchen, H. Decornez, J. R. Furr, and J. Bajorath. Docking and scoring in virtual screening for drug discovery: methods and applications. *Nat Rev Drug Discov*, 3(11):935–949, Nov 2004.
- [73] J. Klokke, P. Langehanenberg, B. Kemper, S. Kosmeier, G. von Bally, C. Riethmuller, F. Wunder, A. Sindic, H. Pavenstadt, E. Schlatter, and B. Edemir. Atrial natriuretic peptide and nitric oxide signaling antagonizes vasopressin-mediated water permeability in inner medullary collecting duct cells. *Am J Physiol Renal Physiol*, 297(3):F693–703, 2009.
- [74] P. Kollman. Free energy calculations: Applications to chemical and biochemical phenomena. *Chem Rev*, 93(7):2395–2417, 1993.
- [75] H. Kubinyi. From narcosis to hyperspace: The history of QSAR. *Quantitative Structure-Activity Relationships*, 21(4):348–356, 2002.
- [76] H. Kubinyi. High throughput in drug discovery. *Drug Discov Today*, 7(13):707–709, Jul 2002.
- [77] S. Kumar, J. M. Rosenberg, D. Bouzida, R. H. Swendsen, and P. A. Kollman. The weighted histogram analysis method for free-energy calculations on biomolecules. I. The method. *J of Comput Chem*, 13:1011–1021, 1992.
- [78] A. Kuo, J. M. Gulbis, J. F. Antcliff, T. Rahman, E. D. Lowe, J. Zimmer, J. Cuthbertson, F. M. Ashcroft, T. Ezaki, and D. A. Doyle. Crystal structure of the potassium channel KirBac1.1 in the closed state. *Science*, 300(5627):1922–1926, Jun 2003.
- [79] K. Lindorff-Larsen, P. Maragakis, S. Piana, M. P. Eastwood, R. O. Dror, and D. E. Shaw. Systematic validation of protein force fields against experimental data. *PLoS One*, 7(2):e32131, 2012.
- [80] K. Lindorff-Larsen, S. Piana, K. Palmo, P. Maragakis, J. L. Klepeis, R. O. Dror, and D. E. Shaw. Improved side-chain torsion potentials for the Amber ff99SB protein force field. *Proteins*, 78(8):1950–8, 2010.
- [81] C. A. Lipinski, F. Lombardo, B. W. Dominy, and P. J. Feeney. Experimental and computational approaches to estimate solubility and permeability in drug discovery and development settings. *Adv Drug Deliv Rev*, 46(1-3):3–26, Mar 2001.
- [82] H. Liu, Z.-B. Gao, Z. Yao, S. Zheng, Y. Li, W. Zhu, X. Tan, X. Luo, J. Shen, K. Chen, G.-Y. Hu, and H. Jiang. Discovering potassium channel

- blockers from synthetic compound database by using structure-based virtual screening in conjunction with electrophysiological assay. *J Med Chem*, 50(1):83–93, Jan 2007.
- [83] H. Liu, Y. Li, M. Song, X. Tan, F. Cheng, S. Zheng, J. Shen, X. Luo, R. Ji, J. Yue, G. Hu, H. Jiang, and K. Chen. Structure-based discovery of potassium channel blockers from natural products: virtual screening and electrophysiological assay testing. *Chem Biol*, 10(11):1103–1113, Nov 2003.
- [84] S. B. Long, E. B. Campbell, and R. Mackinnon. Crystal structure of a mammalian voltage-dependent Shaker family K⁺ channel. *Science*, 309(5736):897–903, Aug 2005.
- [85] S. B. Long, X. Tao, E. B. Campbell, and R. MacKinnon. Atomic structure of a voltage-dependent K⁺ channel in a lipid membrane-like environment. *Nature*, 450(7168):376–382, Nov 2007.
- [86] V. B. Luzhkov and J. Aqvist. Mechanisms of tetraethylammonium ion block in the KcsA potassium channel. *FEBS Lett*, 495(3):191–196, Apr 2001.
- [87] P. D. Lyne, P. W. Kenny, D. A. Cosgrove, C. Deng, S. Zabudoff, J. J. Wendoloski, and S. Ashwell. Identification of compounds with nanomolar binding affinity for checkpoint kinase-1 using knowledge-based virtual screening. *J Med Chem*, 47(8):1962–1968, Apr 2004.
- [88] B. Ma, Y. Xiang, S. M. Mu, T. Li, H. M. Yu, and X. J. Li. Effects of acetazolamide and anordiol on osmotic water permeability in AQP1-cRNA injected *Xenopus* oocyte. *Acta Pharmacol Sin*, 25(1):90–7, 2004.
- [89] J. L. Medina-Franco, K. Martinez-Mayorga, M. A. Giulianotti, R. A. Houghten, and C. Pinilla. Visualization of the chemical space in drug discovery. *Current Computer - Aided Drug Design*, 4(4):322–333, 2008.
- [90] E. C. Meng, B. K. Shoichet, and I. D. Kuntz. Automated docking with grid-based energy evaluation. *Journal of Computational Chemistry*, 13(4):505–524, 1992.
- [91] E. Migliati, N. Meurice, P. DuBois, J. S. Fang, S. Somasekharan, E. Beckett, G. Flynn, and A. J. Yool. Inhibition of aquaporin-1 and aquaporin-4 water permeability by a derivative of the loop diuretic bumetanide acting at an internal pore-occluding binding site. *Mol Pharmacol*, 76(1):105–12, 2009.

- [92] S. Miyamoto and P. A. Kollman. Settle: An analytical version of the SHAKE and RATTLE algorithm for rigid water models. *J Comp. Chem.*, 13:952–962, 1992.
- [93] D. L. Mobley and K. A. Dill. Binding of small-molecule ligands to proteins: "what you see" is not always "what you get". *Structure*, 17(4):489–498, Apr 2009.
- [94] N. Moitessier, P. Englebienne, D. Lee, J. Lawandi, and C. R. Corbeil. Towards the development of universal, fast and highly accurate docking/scoring methods: a long way to go. *Br J Pharmacol*, 153 Suppl 1:S7–26, Mar 2008.
- [95] M. G. Mola, G. P. Nicchia, M. Svelto, D. C. Spray, and A. Frigeri. Automated cell-based assay for screening of aquaporin inhibitors. *Anal Chem*, 81:8219–8229, 2009.
- [96] D. Nelson and M. Cox. *Lehninger Biochemie*. Springer, 4. edition, 2009.
- [97] A. Oda, K. Tsuchida, T. Takakura, N. Yamaotsu, and S. Hirono. Comparison of consensus scoring strategies for evaluating computational models of protein-ligand complexes. *J Chem Inf Model*, 46(1):380–391, 2006.
- [98] M. Parrinello and A. Rahman. Crystal structure and pair potentials: A molecular-dynamics study. *Phys Rev Lett*, 45:1196–1199, 1980.
- [99] M. Parrinello and A. Rahman. Polymorphic transitions in single crystals: A new molecular dynamics method. *J Appl Phys*, 52:7182–7191, 1981.
- [100] B. G. Pierce, Y. Hourai, and Z. Weng. Accelerating protein docking in ZDOCK using an advanced 3D convolution library. *PLoS One*, 6(9):e24657, 2011. zdock.
- [101] M. Rarey, B. Kramer, T. Lengauer, and G. Klebe. A fast flexible docking method using an incremental construction algorithm. *J Mol Biol*, 261(3):470–489, 1996.
- [102] M. N. Rasband, J. S. Trimmer, E. Peles, S. R. Levinson, and P. Shrager. K⁺ channel distribution and clustering in developing and hypomyelinated axons of the optic nerve. *J Neurocytol*, 28(4-5):319–331, 1999.
- [103] K. J. Rhodes, B. W. Strassle, M. M. Monaghan, Z. Bekele-Arcuri, M. F. Matos, and J. S. Trimmer. Association and colocalization of the Kvbeta1

- and Kvbeta2 beta-subunits with Kv1 alpha-subunits in mammalian brain K⁺ channel complexes. *J Neurosci*, 17(21):8246–8258, Nov 1997.
- [104] D. J. Rogers and T. T. Tanimoto. A computer program for classifying plants. *Science*, 132:1115–1118, 1960.
- [105] A. M. Rojek, M. T. Skowronski, E. M. Fuchtbauer, A. C. Fuchtbauer, R. A. Fenton, P. Agre, J. Frokiaer, and S. Nielsen. Defective glycerol metabolism in aquaporin 9 (aqp9) knockout mice. *Proc Natl Acad Sci U S A*, 104(9):3609–14, 2007.
- [106] J. Sadowski, J. Gasteiger, and G. Klebe. Comparison of automatic three-dimensional model builders using 639 X-ray structures. *Journal of Chemical Information and Computer Sciences*, 34(4):1000–1008, 1994.
- [107] J. K. G. Sadowski, J.; Gasteiger. Comparison of automatic three-dimensional model builders using 639 X-ray structures. *J Chem Inf Comput Sci*, 34:1000–1008, 1994.
- [108] M. S. P. Sansom, I. H. Shrivastava, J. N. Bright, J. Tate, C. E. Capener, and P. C. Biggin. Potassium channels: structures, models, simulations. *Biochim Biophys Acta*, 1565(2):294–307, Oct 2002.
- [109] S. S. Scherer and E. J. Arroyo. Recent progress on the molecular organization of myelinated axons. *J Peripher Nerv Syst*, 7(1):1–12, Mar 2002.
- [110] N. Schneider, G. Lange, R. Klein, C. Lemmen, and M. Rarey. HYDEing the false positives-scoring for lead optimization. *J Cheminformatics*, 3:29, 2010.
- [111] G. K. Schuurman-Wolters and B. Poolman. Substrate specificity and ionic regulation of GlnPQ from *Lactococcus lactis*. An ATP-binding cassette transporter with four extracytoplasmic substrate-binding domains. *J Biol Chem*, 280(25):23785–23790, Jun 2005.
- [112] G. K. Schuurman-Wolters, A. Vujičić-Žagar, D.-J. Slotboom, and B. Poolman. Function and structure of the tandem extracytoplasmic substrate-binding domains of the ABC transporter GlnPQ. (in preparation).
- [113] D. Seeliger and B. L. de Groot. Ligand docking and binding site analysis with PyMOL and Autodock-Vina. *J Comput Aided Mol Des*, 24(5):417–422, May 2010.

- [114] Y. Shan, E. T. Kim, M. P. Eastwood, R. O. Dror, M. A. Seeliger, and D. E. Shaw. How does a drug molecule find its target binding site? *J Am Chem Soc*, 133(24):9181–9183, 2011.
- [115] B. K. Shoichet. Virtual screening of chemical libraries. *Nature*, 432(7019):862–865, Dec 2004.
- [116] M. A. Sills, D. Weiss, Q. Pham, R. Schweitzer, X. Wu, and J. J. Wu. Comparison of assay technologies for a tyrosine kinase assay generates different results in high throughput screening. *J Biomol Screen*, 7(3):191–214, Jun 2002.
- [117] M. Stahl and M. Rarey. Detailed analysis of scoring functions for virtual screening. *J Med Chem*, 44(7):1035–1042, Mar 2001.
- [118] G. S. Tamura, A. Nittayajarn, and D. L. Schoentag. A glutamine transport gene, *glnQ*, is required for fibronectin adherence and virulence of group b *Streptococci*. *Infection and Immunity*, 70(6):2877–2885, 2002.
- [119] P. D. Thomas and K. A. Dill. Statistical potentials extracted from protein structures: how accurate are they? *J Mol Biol*, 257(2):457–469, Mar 1996.
- [120] G. Torrie and J. Valleau. Nonphysical sampling distributions in monte carlo free-energy estimation: Umbrella sampling. *J Comput Phys*, 23(2):187 – 199, 1977.
- [121] O. Trott and A. J. Olson. AutoDock Vina: improving the speed and accuracy of docking with a new scoring function, efficient optimization, and multithreading. *J Comput Chem*, 31(2):455–461, Jan 2010.
- [122] H. Tsukaguchi, C. Shayakul, U. V. Berger, B. Mackenzie, S. Devidas, W. B. Guggino, A. N. van Hoek, and M. A. Hediger. Molecular characterization of a broad selectivity neutral solute channel. *J Biol Chem*, 273(38):24737–24743, Sep 1998.
- [123] T. van der Heide and B. Poolman. ABC transporters: one, two or four extracytoplasmic substrate-binding sites? *EMBO Rep*, 3(10):938–943, Oct 2002.
- [124] D. Van Der Spoel, E. Lindahl, B. Hess, G. Groenhof, A. Mark, and H. Berendsen. GROMACS: fast, flexible, and free. *Journal of computational chemistry*, 26(16):1701, 2005.

- [125] W. van Gunsteren and H. Berendsen. Computer simulation of molecular dynamics: Methodology, applications, and perspectives in chemistry. *Angew Chem Int Ed Engl*, 29:992–1023, 1990.
- [126] H. J. Verheij. Leadlikeness and structural diversity of synthetic screening libraries. *Mol Divers*, 10(3):377–388, Aug 2006.
- [127] A. S. Verkman. Drug discovery in academia. *Am J Physiol Cell Physiol*, 286(3):C465–C474, Mar 2004.
- [128] S. J. Wacker, W. Jurkowski, K. J. Simmons, C. W. G. Fishwick, A. P. Johnson, D. Madge, E. Lindahl, J.-F. Rolland, and B. L. de Groot. Identification of selective inhibitors of the potassium channel Kv1.1-1.2((3)) by high-throughput virtual screening and automated patch clamp. *ChemMed-Chem*, -(in print), Mar 2012.
- [129] T. Walz, Y. Fujiyoshi, and A. Engel. The AQP structure and functional implications. *Handb Exp Pharmacol*, 190(190):31–56, 2009.
- [130] H. Wang, D. D. Kunkel, T. M. Martin, P. A. Schwartzkroin, and B. L. Tempel. Heteromultimeric K⁺ channels in terminal and juxtaparanodal regions of neurons. *Nature*, 365(6441):75–79, Sep 1993.
- [131] H. Wang, D. D. Kunkel, P. A. Schwartzkroin, and B. L. Tempel. Localization of Kv1.1 and Kv1.2, two K channel proteins, to synaptic terminals, somata, and dendrites in the mouse brain. *J Neurosci*, 14(8):4588–4599, Aug 1994.
- [132] J. Wang, P. Cieplak, and P. A. Kollman. How well does a restrained electrostatic potential (RESP) model perform in calculating conformational energies of organic and biological molecules? *J comput chem*, 21:1049 – 1074, 2000.
- [133] J. Wang, W. Wang, P. A. Kollman, and D. A. Case. Automatic atom type and bond type perception in molecular mechanical calculations. *J Mol Graph Model*, 25(2):247–60, 2006.
- [134] J. Wang, R. M. Wolf, J. W. Caldwell, P. A. Kollman, and D. A. Case. Development and testing of a general amber force field. *J Comput Chem*, 25(9):1157–1174, Jul 2004.

- [135] R. Wang and S. Wang. How does consensus scoring work for virtual library screening? an idealized computer experiment. *J Chem Inf Comput Sci*, 41(5):1422–1426, 2001.
- [136] E. Winsberg. *Science in the Age of Computer Simulation*. The University of Chicago Press, 2010.
- [137] J.-M. Yang, Y.-F. Chen, T.-W. Shen, B. S. Kristal, and D. F. Hsu. Consensus scoring criteria for improving enrichment in virtual screening. *Journal of Chemical Information and Modeling*, 45(4):1134–1146, 2005.
- [138] Q. Yang, L. Du, X. Wang, M. Li, and Q. You. Modeling the binding modes of Kv1.5 potassium channel and blockers. *J Mol Graph Model*, 27(2):178–187, Sep 2008.
- [139] Q. Yang, D. Fedida, H. Xu, B. Wang, L. Du, X. Wang, M. Li, and Q. You. Structure-based virtual screening and electrophysiological evaluation of new chemotypes of K(v)1.5 channel blockers. *ChemMedChem*, 5(8):1353–1358, Aug 2010.
- [140] H.-X. Zhou and M. K. Gilson. Theory of free energy and entropy in non-covalent binding. *Chem Rev*, 109(9):4092–4107, Sep 2009.
- [141] Y. Zhou, J. H. Morais-Cabral, A. Kaufman, and R. MacKinnon. Chemistry of ion coordination and hydration revealed by a K⁺ channel-Fab complex at 2.0 Å resolution. *Nature*, 414(6859):43–48, Nov 2001.
- [142] F. Zhu, E. Tajkhorshid, and K. Schulten. Collective diffusion model for water permeation through microscopic channels. *Phys Rev Lett*, 93(22):224501, Nov 2004.
- [143] M. Zink and H. Grubmüller. Mechanical properties of the icosahedral shell of southern bean mosaic virus: a molecular dynamics study. *Biophys J*, 96(4):1350–1363, Feb 2009.
- [144] M. Zink and H. Grubmüller. Primary changes of the mechanical properties of Southern Bean Mosaic Virus upon calcium removal. *Biophys J*, 98(4):687–695, Feb 2010.
- [145] M. Zolli-Juran, J. D. Cechetto, R. Hartlen, D. M. Daigle, and E. D. Brown. High throughput screening identifies novel inhibitors of Escherichia coli dihydrofolate reductase that are competitive with dihydrofolate. *Bioorg Med Chem Lett*, 13(15):2493–2496, Aug 2003.

- [146] Z. Zsoldos, D. Reid, A. Simon, B. S. Sadjad, and A. P. Johnson. eHiTS: an innovative approach to the docking and scoring function problems. *Curr Protein Pept Sci*, 7(5):421–435, Oct 2006.
- [147] Z. Zsoldos, D. Reid, A. Simon, S. B. Sadjad, and A. P. Johnson. eHiTS: a new fast, exhaustive flexible ligand docking system. *J Mol Graph Model*, 26(1):198–212, Jul 2007.
- [148] R. W. Zwanzig. High-temperature equation of state by a perturbation method. I. Nonpolar gases. *J Chem Phys*, 22:1420, 1954.

Danksagung

An dieser Stelle möchte ich all jenen danken, die zum erfolgreichen Abschluß dieser Arbeit beigetragen haben. An erster und herausgehobener Stelle möchte ich meinen Eltern und meiner Schwester danken. Eure jahrelange Unterstützung in all meinen Entscheidungen gab mir Sicherheit und Motivation. Danke, dass ich mich immer auf euch verlassen kann.

Bei Dr. Bert de Groot möchte ich mich ganz besonders herzlich bedanken. Er stand mir in den letzten Jahren mit fruchtbaren Ratschlägen und Ideen zur Seite und hatte stets eine offene Tür und Zeit für Diskussionen. Die in seiner Arbeitsgruppe gebotene Freiheit, kann ich jedem engagierten Doktoranden nur wünschen. Ich hatte immer das Gefühl bei den verschiedenen Projekten kompetent und vorausschauend beraten zu sein, um in eigener Verantwortung die nächsten Etappen auswählen und entwickeln zu können. Eine bessere Umgebung für innovative Forschung kann ich mir schwer vorstellen.

Bei Dr. Helmut Grubmüller möchte ich mich für die einzigartige Atmosphäre in seiner Abteilung bedanken. Ebenso danke ich den Professoren meines Prüfungsausschusses Dr. Jörg Enderlein und Dr. Holger Stark für ihre Zeit und den Aufwand welche sie für Begutachtung meiner Arbeit aufgewendet haben, sowie für ihre Anregungen und Ideen.

Großer Dank gilt auch unserer Sekretärin Eveline Heinemann, die mir im Verwaltungskrieg mit Formularen, Anträgen und Genehmigungen stets eine helfende Hand und notfalls die rettende Verbündete war. Desweiteren danke ich Ansgar Esztermann und Martin Fechner, unseren Administratoren, welche sich liebevoll um die gut 15.000 Computerkerne in unserer Abteilung kümmern und stets jegwede soft- oder hardware basierte Anforderung technisch realisierten.

Außerdem danke ich den Mitgliedern der Abteilung für theoretische und computergestützte Biophysik, insbesondere meinen Kollegen Dr. Carsten Kutzner und Timo Graen, mit denen ich eines der schönsten (um nicht den Elativ, den absoluten Superlativ, zu gebrauchen) Büros beziehen durfte, und denen ich für die wunderbare Arbeitsatmosphäre ganz besonders danke.

Nicht zuletzt möchte ich den zahlreichen Wissenschaftlern danken, die in professioneller Zusammenarbeit die Verwirklichung dieser Arbeit ermöglicht haben. Dies waren insbesondere Dr. Wiktor Jurkowski, Dr. Katie Simmons und Jean-Francois Rolland in dem Kaliumkanal Projekt, sowie Dr. Michael Rützler und Dawid Krenc, die maßgeblich am Aquaporin 9 Projekt beteiligt waren. Die Experimente im Projekt mit dem ABC-Transporter GlnPQ wurden durchgeführt Dr. Bert Poolman und Faizah Fulyani. Herausstellen möchte ich dabei die enorm produktive Arbeit mit Dr. Michael Rützler, welche mir eine ausserordentliche Freude bereitete. Dr. Daniel Seeliger und Dr. Camilo Aponte-Santamaria stellten technische Unterstützung und Computerskripte bereit. Außerdem möchte ich all jenen danken, die diese Arbeit in den letzten Monaten korrektur gelesen haben, Dr. Petra Kellers, Timo Graen, David Köpfer, Dr. Hadas Leonov, Colin A. Smith, Vytautas Gapsys, Shreyas Kaptan, Dr. Gregory Bubnis und Michal Walczak.

Zu guter Letzt möchte ich hier meinen Freunden in und um Göttingen danken. Die Unterstützung, die ich von eurer Seite erfahren habe, war von einer ganz anderen Qualität und gehört zweifellos zu den wichtigsten Erfahrungen meines Lebens.

Sören Wacker

Göttingen, July 2012

This research was funded by the European Community's Seventh Framework Program FP7/2007-2013 under grand agreement no. HEALTH-F4-2007-201924, EDICT Consortium.

Appendix

Table 7.1: BEDROC and AROC values according to scores and subscores of Vina, Glide, eHiTS and FlexX.

Score/Sub-Score	BEDROC	AROC
Autodock-Vina	0.519	0.706
FlexX-Ambig-Score	0.236	0.559
FlexX-AnzMatch	0.107	0.409
FlexX-AvgVolume	0.133	0.5
FlexX-Clash-Score	0.042	0.392
FlexX-Lipo-Score	0.59	0.703
FlexX-Match-Score	0.023	0.372
FlexX-MaxVolume	0.09	0.449
FlexX-Rot-Score	0.187	0.431
FlexX-Total-Score	0.109	0.439
ChemScore-Clash-Score	0.105	0.464
ChemScore-FragNo	0.239	0.355
ChemScore-Lipo-Score	0.549	0.731
ChemScore-Match-Score	0.011	0.327
ChemScore-MaxVolume	0.107	0.492
ChemScore-Rot-Score	0.219	0.466
ChemScore-Total-Score	0.336	0.632
Glide-XP-Electro	0.026	0.379
Glide-XP-HBond	0.022	0.38
Glide-XP-LipophilicEvdW	0.606	0.725
Glide-XP-Penalties	0.174	0.405
Glide-XP-PhobicPenal	0.206	0.51
Glide-XP-RotPenal	0.07	0.502
Glide-XP-Sitemap	0.358	0.6
Glide-ecoul	0.026	0.385
Glide-einternal	0.11	0.392
Glide-emodel	0.561	0.707
Glide-energy	0.539	0.695
Glide-evdw	0.67	0.741
Glide-gscore	0.315	0.548
eHiTS-Energy	0.336	0.614
eHiTS-Score	0.165	0.504
eHiTS-Term-Coulomb	0.163	0.542
eHiTS-Term-H-bond	0.089	0.432
eHiTS-Term-Lcover	0.131	0.48
eHiTS-Term-Lipophil	0.124	0.454
eHiTS-Term-LlogD	0.086	0.387
eHiTS-Term-Rcharge	0.053	0.372
eHiTS-Term-Rcover	0.23	0.475
eHiTS-Term-RlogD	0.245	0.536
eHiTS-Term-Rshape	0.298	0.563
eHiTS-Term-depth	0.44	0.693
eHiTS-Term-entropy	0.051	0.34
eHiTS-Term-family	0.446	0.703
eHiTS-Term-other	0.212	0.541
eHiTS-Term-pi-stack	0.241	0.533
eHiTS-Term-solvent	0.252	0.606
eHiTS-Term-steric	0.303	0.55
eHiTS-Term-strain	0.582	0.735

Table 7.2: Other known $K_V1.1$ and $K_V1.2$ inhibitors and references.

Name/INN	CAS	Pubmed ID	$K_V1.1$ IC ₅₀	$K_V1.2$ IC ₅₀	Other tar-gets IC ₅₀
Nifedipine	21829-25-4	7517498	NA	NA	NA
Amitriptyline	50-48-6	17456683	22	-	$K_V7.2/7.3$ 10
4-AP	504-24-5	19413590	170	230	$K_V1.4, K_V4.2$
Fampridine	119229-65-1	19413590	3.6	3.7	NavX 11.9
Nerispiridine	178894-81-0	15967421	NA	NA	NA
AM92016 hydrochloride	139298-40-1	16368898	-	1	$K_V1.4,$ $K_V1.5,$ $K_V2.1,$ $K_V3.2,$ $K_V4.2,$ hERG
KN-93					

Table 7.3: Original ranks of the 14 hits according to the implementations A and B and the three used consensus scoring methods, rank2max, rank2number and rank2rank.

ID	Implementation A			Implementation B			IC50
	rank2max	rank2number	rank2rank	rank2max	rank2number	rank2rank	
1	94	216	55	n.e.	n.e.	n.e.	0.71
2	121	373	225	n.e.	n.e.	n.e.	0.79
3	86	575	181	n.e.	n.e.	n.e.	1.41
4	170	228	208	n.e.	n.e.	n.e.	1.62
5	72	455	132	n.e.	n.e.	n.e.	2.98
6	91	410	91	n.e.	n.e.	n.e.	4.07
7	80	583	152	8	35	10	1.53
8	2399	2084	1137	68	184	115	0.58
9	418	1086	333	9	70	23	0.93
10	323	471	315	7	23	8	1.66
11	1204	782	352	133	88	73	2.71
12	3563	1805	1401	137	112	53	3.7
13	267	740	444	13	51	30	3.77
14	303	754	354	112	78	100	5.94

Table 7.4: Measured inhibition of the Maybridge compounds at a given concentration.

Maybridge ID	% Inhibition [%]	At conc. [μ M]
CD05595	30	100
DFP00270	30	100
HTS03850	100	100
HTS06008	50	100
HTS07176	50	100
HTS13286	100	25
phloretin	80	100
RF03176	100	25
XBX00246	50	25
HTS13772	100	100

Table 7.5: Compounds inactive on murine AQP9 included in the set of random decoys.

W00328	inactive on mAQP9
BTB06069	inactive on mAQP9
BTB08755	inactive on mAQP9
GK01664	inactive on mAQP9
GK02411	inactive on mAQP9
HTS01865	inactive on mAQP9
HTS03989	inactive on mAQP9
HTS04489	inactive on mAQP9
HTS05750	inactive on mAQP9
HTS06687	inactive on mAQP9
HTS12517	inactive on mAQP9
KM02905	inactive on mAQP9
KM03611	inactive on mAQP9
KM04183	inactive on mAQP9
KM06426	inactive on mAQP9
KM08508	inactive on mAQP9
KM09574	inactive on mAQP9
KM10521	inactive on mAQP9
ML00240	inactive on mAQP9
RDR03718	inactive on mAQP9
S12781	inactive on mAQP9
S07440	inactive on mAQP9
S07673	inactive on mAQP9
S12781	inactive on mAQP9
SEW00445	inactive on mAQP9
SEW03679	inactive on mAQP9
SP01460	inactive on mAQP9
SPB06953	inactive on mAQP9
SPB08436	inactive on mAQP9
TB00031	inactive on mAQP9

Table 7.6: Screening results at 100 μ M concentration (part I).

Enamine ID	Half Life 1	Half Life 2	Half Life 3	ANOVA and Dunnet's test	Average	Rank
T6963384	13.1	9.164	10.12	***	10.795	1
T6666911	8.031	9.061	8.432	***	8.508	2
T6616714	7.43	6.059	10.49	***	7.993	3
T6473159	6.61	6.827	7.948	***	7.128	4
T6945453	6.783	5.367	6.97	***	6.373	5
T6406844	5.781	5.139	6.182	***	5.701	6
T6442514	4.838	4.301	5.2	***	4.78	7
T6655310	4.476	3.506	6.067	***	4.683	8
T6753522	4.683	4.975	3.445	***	4.368	9
T6674497	4.635	2.868	4.165	***	3.889	10
T6374888	3.578	3.337	3.338	***	3.418	11
T6617912	3.215	2.606	4.16	***	3.327	12
T6912201	3.727	2.969	3.015	***	3.237	13
T6837949	3.703	2.537	3.44	***	3.227	14
T6090579	3.151	3.008	2.519	***	2.893	15
T6878499	3.03	2.318	3.268	**	2.872	16
T6662346	2.637	3.395	1.777	***	2.603	17
T6836856	3.301	1.904	2.464	***	2.556	18
T6797755	3.367	2.564	1.498	*	2.476	19
T6792287	2.527	2.016	2.351	ns	2.298	20
T6782182	2.382	2.65	1.796	ns	2.276	21
T6669107	2.234	2.17	2.183	**	2.196	22
T6237899	1.9	1.845	1.82	***	1.855	23
T6929930	2.309	1.622	1.623	**	1.851	24
T6453621	1.89	2.033	1.55	***	1.824	25
T6344977	1.521	2.067	1.785	***	1.791	26
T6470899	2.415	1.447	1.498	**	1.787	27
T6424708	1.599	1.717	1.942	***	1.753	28
T6280980	2.009	1.424	1.659	ns	1.697	29
T5925407	1.853	1.369	1.86	ns	1.694	30
T5819367	2.262	1.208	1.597	*	1.689	31
T6229572	1.67	1.23	2.003	ns	1.634	32
T6848089	1.734	0.985	1.573	***	1.431	33
T6770322	1.301	1.322	1.361	**	1.328	34
T6217877	1.467	1.045	1.397	**	1.303	35
T6593327	1.276	1.177	1.261	ns	1.238	36
T6130871	1.377	1.005	1.208	ns	1.197	37
T6208562	1.215	0.954	1.373	ns	1.181	38
T6643668	1.114	0.852	1.309	ns	1.092	39
T6661714	1.103	0.87	1.251	ns	1.075	40
T6706567	1.143	0.883	1.172	ns	1.066	41
T6778132	1.062	0.837	1.218	ns	1.039	42
T6591224	1.322	0.981	0.803	ns	1.035	43
T6602468	1.267	0.873	0.92	ns	1.02	44
T6299997	1.196	0.787	1.033	ns	1.005	45
T6063556	1.094	0.885	0.984	ns	0.987	46
T6604432	1.131	0.823	0.991	ns	0.981	47
T6389731	1.387	0.74	0.813	ns	0.98	48
T6718086	0.803	1.294	0.791	ns	0.963	49
T6496774	1.02	0.659	1.156	ns	0.945	50
T6648591	1.011	0.717	1.066	ns	0.931	51
T6244162	1.001	0.882	0.907	ns	0.93	52
T5290442	0.889	1.055	0.826	ns	0.923	53
T6792868	1.126	0.773	0.852	ns	0.917	54
T6506095	0.922	0.775	1.045	ns	0.914	55

Table 7.7: Screening results at 100 μ M concentration (part II).

Enamine ID	Half Life 1	Half Life 2	Half Life 3	ANOVA and Dunnet's test	Average	Rank
T6963384	13.1	9.164	10.12	***	10.795	1
T6169898	0.928	0.767	1.01	ns	0.902	56
T6045127	1.17	0.758	0.756	ns	0.895	57
T5625571	1.099	0.655	0.92	ns	0.892	58
T5564209	0.913	0.798	0.962	ns	0.891	59
T6899185	1.062	0.757	0.802	ns	0.874	60
T5925215	1.061	0.607	0.947	ns	0.872	61
T6551048	0.843	0.876	0.865	ns	0.862	62
T6869847	0.978	0.605	0.932	ns	0.838	63
T6348177	0.901	0.725	0.887	ns	0.838	64
T5935261	1.13	0.569	0.806	ns	0.835	65
T6472369	0.947	0.631	0.921	ns	0.833	66
T6260743	0.917	0.733	0.779	ns	0.81	67
T6400536	0.919	0.69	0.805	ns	0.805	68
T5838470	1.048	0.576	0.78	ns	0.801	69
T5976868	0.838	0.554	1.009	ns	0.8	70
T6389727	0.953	0.667	0.754	ns	0.791	71
T6762259	0.811	0.773	0.745	ns	0.776	72
T6324174	0.807	0.512	0.965	ns	0.761	73
T5273571	0.869	0.527	0.884	ns	0.76	74
T5283260	0.932	0.468	0.869	ns	0.756	75
T6700376	0.786	0.556	0.921	ns	0.754	76
T6718742	0.821	0.73	0.679	ns	0.743	77
T6467923	0.899	0.728	0.601	ns	0.742	78
T6487759	0.867	0.627	0.724	ns	0.739	79
T6655300	0.782	0.575	0.841	ns	0.732	80
T6821126	0.904	0.59	0.685	ns	0.726	81
T6798120	0.885	0.505	0.769	ns	0.72	82
T6502363	0.695	0.491	0.94	ns	0.709	83
T6941479	0.821	0.52	0.756	ns	0.699	84
T5728907	0.843	0.641	0.612	ns	0.698	85
T6696797	0.686	0.676	0.723	ns	0.695	86
T6805975	0.749	0.476	0.852	ns	0.692	87
DMSO	0.715	0.61	0.744	0.69	88	
T6446477	0.87	0.428	0.771	ns	0.69	89
T6204618	0.978	0.472	0.61	ns	0.687	90
T5618514	0.636	0.641	0.754	ns	0.677	91
T6926933	0.779	0.583	0.662	ns	0.675	92
T6604321	0.661	0.566	0.793	ns	0.673	93
T6435541	0.621	0.512	0.856	ns	0.663	94
T6619299	0.534	0.877	0.569	ns	0.66	95
T6309181	0.645	0.48	0.849	ns	0.658	96
T6660708	0.675	0.594	0.659	ns	0.643	97
T6278016	0.666	0.569	0.673	ns	0.636	98
T5765412	0.679	0.519	0.705	ns	0.634	99
T6795119	0.695	0.551	0.657	ns	0.634	100
T6666617	0.735	0.501	0.652	ns	0.63	101
T6421733	0.576	0.545	0.689	ns	0.603	102
T6797784	0.761	0.448	0.569	ns	0.593	103
T5275225	0.577	0.429	0.72	ns	0.575	104
T5392604	0.645	0.513	0.541	ns	0.566	105
T6769024	0.512	0.539	0.527	ns	0.526	106

Table 7.8: Dose response results between 0 and 100 μM concentration.

Enamine ID	IC ₅₀
T6090579	$3.1 \cdot 10^{-06}$
T6674497	$3.93 \cdot 10^{-06}$
T6666911	$7.67 \cdot 10^{-06}$
T6963384	$1.01 \cdot 10^{-05}$
T6616714	$1.44 \cdot 10^{-05}$
T6473159	$1.77 \cdot 10^{-05}$
T6374888	$2.43 \cdot 10^{-05}$
T6878499	$3.83 \cdot 10^{-05}$
T6837949	$4.01 \cdot 10^{-05}$
T6753522	$5.50 \cdot 10^{-05}$
T6617912	$7.37 \cdot 10^{-05}$
T6406844	0.0001499
T6655310	0.0002472
T6836856	0.0007398
T6797755	0.001433
T6912201	0.05073
T6442514	0.05777
T6945453	1.332

Lebenslauf von Sören Wacker

Persönliche Daten

Name: Sören Wacker
Adresse: Untere-Masch-Str. 7, 37073 Göttingen
Geburtsdatum: 22. Januar 1982
Geburtsort: Stadthagen
Nationalität: Deutsch

Ausbildung

01.2009 – 06.2012 Doktorarbeit mit dem Titel:
Computer-Aided Drug Design for Membrane Channel Proteins
Am Max Planck Institut für biophysikalische Chemie in Göttingen
Betreuer: Prof. Dr. Bert de Groot

12.2007 – 12.2008 Diplomarbeit mit dem Titel:
Molecular Dynamics Study of Potassium Channels
An der Georg-August-Universität Göttingen
Betreuer: Prof. Dr. Bert de Groot

11.2006 – 12.2008 Studium der Physik an der Georg-August-Universität Göttingen

11.2002 – 05.2006 Studium der Physik an der Leibniz-Universität Hannover

06.2001 Abitur am Wilhelm-Busch-Gymnasium, Stadthagen

07.1998 Erweiterter Sekundarabschluss I
an der IGS Schaumburg, Stadthagen

