

1 Introduction

1.1 Focal adhesion kinase

Focal adhesions (FA) are macromolecular protein complexes which act as a connection hub between the cell, i.e. the cytoskeleton, and the extracellular matrix (ECM). They enable the cell to performing tension forces, but can also trigger mechanical stimuli from the ECM. One important protein associated to FA is the focal adhesion kinase (FAK). FAK occurs in several signalling pathways and is a key player in integrating extracellular stimuli. It is of large interest not least because in cancer cells often an overexpression of FAK can be found and understanding the activation processes and dynamics of FAK could give rise to new cancer treatments.

1.1.1 Structure

FAK consists of four major domains (see ??): a FERM domain as N-terminal, a tyrosine kinase, a proline rich region and a focal adhesion targeting (FAT) domain as C-terminal.

FERM (4.1 protein, ezrin, radixin and moesin) is a common protein domain, which targets proteins to membranes [9] and consists of three subdomains F1, F2 and F3. In the F2 subdomain there is the basic patch (²¹⁶KAKTLRK²²²), which is a prominent binding site for phosphatidylinositol-4,5-bisphosphate (PI(4,5)P₂). This phospholipid is locally generated in FA due to integrin signaling [15].

The kinase domain of the C-lobe, the activation loop and the N-lobe. The catalytic activity of kinase is regulated by the phosphorylation of Y⁵⁷⁶ and Y⁵⁷⁷, which are located in the activation loop [7]. The C-lobe also provides a binding site for PI(4,5)P₂ which is located next to the basic patch of the FERM domain [12].

The FERM domain and the kinase are connected by a linker region. In contrast to other kinase domains the main autophosphorylation site of FAK Y³⁹⁷ can be found in this region [15].

The FAT domain is linked to the kinase by a flexible proline rich region. FAT targets to FA by interacting with talin and paxillin, which are proteins associated with FA [2].

1.1.2 Autophosphorylation and activation

A maximum catalytic turnover requires Y^{576} and Y^{577} in the activation loop to be phosphorylated. In the inactive state this region is shielded by the FERM domain. Also the autophosphorylation site Y^{397} is isolated by the FERM domain. Therefore an activation is only possible if the FERM domain dissociate at least partly from the kinase [24]. FAK triggers several stimuli, but in this thesis the main focus lies on the allosteric effect of $PI(4,5)P_2$ binding to the basic patch.

$PI(4,5)P_2$ has a net charge of -4, but in presence of K the deprotonated state gets promoted resulting in a net charge of -5. The electrostatic binding of $PI(4,5)P_2$ to the basic patch in the F2 subdomain results in long ranged configurational changes, which also influence the interface between the F1 subdomain and the N-lobe. Also the linker region gets less strongly bound so that an autophosphorylation of Y^{397} is promoted. Phosphorylated Y^{397} is a suitable binding site for SH2, a subdomain of proteins from the Src kinase family. The presence of a kinase and the partial opening of the FERM kinase interface leads to a phosphorylation of Y^{576} and Y^{577} . The resulting Src-FAK complex can act as an fully active kinase. In this state the FERM domain dissociates from the kinase [35, 15].

1.1.3 Dimerization and clustering

The FERM domain induces a dimerization of FAK as it is doing in other proteins containing a FERM domain as well. The interaction emerges around W^{266} in the connected domains and is stabilised by an interaction of the FAT domain with the basic patch of the other FERM domain respectively [6]. It has been shown, that the autophosphorylation happens in trans [33] and requires W^{266} [6]. On a membrane however the dimer can not be stabilised by the FAT-FERM interaction, because the basic patch binds to $PI(4,5)P_2$ in the membrane, so the membrane has to stabilize the dimer. Besides dimers larger clusters of FAK can emerge, which can contribute to additional signaling processes [21].

1.2 Previous work and motivation

1.3 Molecular Dynamics simulations

It is often difficult to get a precise insight into dynamics of a biological system on a molecular scale. Reasons are, that the measurable time- and length scales are usually larger than those of the system. Therefore most effects can only be measured indirectly. Another problem are disturbing factors, which can be hard to control.

These problems do not exist in computer simulations, in which the positions of all particles of the system are known at a much smaller timescale than the dynamics in the system.

An important tool is Molecular dynamics simulation (MD). MD can provide atomistic details, but e.g. with coarse grained models (see [section 1.3.2](#)) a length scale of several nanometres is still accessible for microseconds. Therefore MD is suitable for membranes, proteins and other biomolecules.

In this thesis MD was used to investigate interactions of FAK with other FAK and the membrane. Therefore the main concepts of MD and characteristics of the used models are outlined below. However within the scope of this thesis it can just be a small survey.

Because GROMACS (Groningen Machine for Chemical Simulations) [[1](#), [3](#)] was used as MD engine in this thesis, the following explanations refer to GROMACS conventions and features.

1.3.1 The physics behind MD

In MD the system is simulated on an atomistic scale but with classical mechanics only. In order to do so atoms are only treated as spheres without distinction into electrons and nuclei. Quantum mechanical (QM) effects, such as excitation or electron transfer processes, are therefore not accessible and neglected in MD. However the atoms are parametrized by effective parameters, which can be motivated from QM [[28](#), p. 127f].

Newtons equation of motion

$$\vec{F}_i = m_i \frac{d^2 \vec{r}_i}{dt^2} \quad (1.1)$$

where \vec{F}_i is the force acting on particle i , m_i its mass and \vec{r}_i its position, can be turned into two first-order differential equations

$$\frac{d\vec{r}_i}{dt} = \vec{v}_i \quad (1.2)$$

$$m_i \frac{d\vec{v}_i}{dt} = \vec{F}_i \quad (1.3)$$

They can be integrated numerically with e.g. Leapfrog- or Verlet integration scheme. Both are low order, i.e. the integration error is $\mathcal{O}(\Delta t^3)$, which saves computational cost. However they have the great advantage to be time reversible and symplectic. Later means, that it conserves the phase space volume, like the Hamiltonian operator would do as well. Therefore the long term error in energy conservation stays small [14, p. 72ff].

The force acting on particle i is given by the potential at its position.

$$\vec{F}_i = -\frac{\partial V}{\partial \vec{r}_i} \quad (1.4)$$

In MD bonded and non-bonded interactions contribute to the potential V , but it is also possible to apply external forces.

Bonded interactions

Bonded interactions act intra molecular and describe chemical bonds. They can occur between two, three or four particles.

An interaction between two bonded particles refers to their bond length. A deviation from the equilibrium bond length results in potential energy, which is usually described by an harmonic oscillator.

$$V_{\text{dist, bond } i} = \frac{k_{\text{dist}}}{2} (r - r_0)^2 \quad (1.5)$$

k_{dist} is the force constant and r_0 is the equilibrium bond length for the bond type of bond i . For larger deviations a Morse potential (exponential decaying potential for large deviations) is more precise, but has a much higher computational cost.

Also a deviation from an equilibrium angle between three bonded partners, i.e. the bond angle, results in potential energy. A common description is the harmonic

oscillator as well.

$$V_{\text{angle}} = \frac{k_{\text{angle}}}{2} (\theta - \theta_0)^2 \quad (1.6)$$

The dihedral angle is the angle between two particles, which are separated by three bonds and can therefore be understood as the torsion angle of the intermediate two particles and bond. Besides a description of torsion this angle can also be used to preserve plane rings and the chirality of four particle groups. It is usually approximated with a periodic approach

$$V_{\text{dihedral, periodic}} = \frac{k_{\text{dihedral}}}{2} (1 + \cos(n\phi - \phi_0)) \quad (1.7)$$

where k_{dihedral} describes the energy barrier for turning the dihedral angle, n the number of minima in the energy function (multiplicity) and ϕ_0 a phase factor [1, p. 71-83].

Non-bonded interactions

Non-bonded interactions are present between all atoms in the system and act pairwise. In MD Pauli repulsion, van der Waals (vdW) forces and electrostatic forces are taken into account. Because bonded interactions use effective parameters they are usually excluded from non-bonded interactions.

The Lennard-Jones potential combines Pauli repulsion (r^{-12} term) and the vdW force (r^{-6} term).

$$V_{\text{Lennard-Jones}} = \sum_{\text{non-bonded pairs } i,j} 4\epsilon \left(\left(\frac{\sigma}{r_{ij}} \right)^{12} - \left(\frac{\sigma}{r_{ij}} \right)^6 \right) \quad (1.8)$$

ϵ is related to the potential depth and σ to the potential range.

The Coulomb potential is given by

$$V_{\text{Coulomb}} = \frac{q_1 q_2}{4\pi\epsilon_0\epsilon_r r} \quad (1.9)$$

where q_1 , q_2 are the (partial) charges of the interacting particles, r their distance and ϵ_r the relative dielectric constant [1, p. 65-71].

In general non-bonded interactions act between all atoms in the system, which come along with a very large computational cost.

The easiest solution is to use a cut-off radius r_c . Particles behind this radius are not taken into account. This can be implemented very efficiently with Verlet neighbour lists. For each particle a neighbour list is created, which contains all particles inside a second radius r_v with $r_v > r_c$. For a force calculation only the distances to the particles, which are part of the list, have to be calculated. The lists are updated, if the maximal displacement in the system is larger than $r_v - r_c$. This method is suitable especially for Lennard-Jones potential, because of its rapid decay r_c can be chosen very small [28, p. 144].

The electrostatic potential is proportional to $1/r$, that is why the use of a cut-off radius would lead to large jumps in the potential. Long range interactions have to be considered, which can be effectively done with Particle Mesh Ewald (PME) summation [10]. Particle mesh methods in general split the electrostatic potential up into a short range and a long range part via a switching function. The short range part can be calculated with a small cut-off radius in real space. The long range part however is calculated by solving the Poisson equation of the actual charge distribution, for which a discrete grid (mesh) is used. In PME this grid is transformed to Fourier space, where the solution of the Poisson equation is a sum over the gridpoints. This requires of course periodic boundary conditions (see section 1.3.3) [28, p. 246-251].

External forces

In addition to bonded and non-bonded interactions external or artificial forces can be applied, such as pulling forces, restrains and constrains.

With GROMACS it is possible to perform pulling forces onto groups of atoms in the system. In this thesis pulling was used to bias distances between groups. For this GROMACS provides an option to apply an umbrella potential to two groups which yields in a force, which is proportional to the deviation of the distance between the groups from a reference distance. The force can be applied on one or two spatial dimensions only or along a predefined vector and the reference distance can change during in time [1, p. 154-159].

Restraints are artificial potentials applied to positions, distances, angles, dihedral angles or orientations of particle groups to conserve labile configurations or to take additional information from experimental data into account [1, p. 84f].

Constraints are set to keep properties, such as bond lengths constant. In order to do so, they get reset to the desired value after each time step [1, p. 44f].

1.3.2 Forcefields

The parameters for the potentials described above are provided by force fields. There is a wide range of force fields, which are optimized for different application fields. In this thesis two different ones, CHARMM36 and Martini, are introduced.

Force fields define specific atomtypes, which allow a mapping of atoms (particles in physical system) onto beads (particles used in the simulation). This mapping can take the environment of the particle (binding partners, solvents, nearby charges a.o.) into account, but can also neglect details by f.e. mapping several atoms to one bead (coarse-graining).

Force fields not only define the atomtypes and their properties, but also all parameters for the calculation of the potential, especially force constants, equilibrium distances a.o. Therefore they define the physics of the system.

All-atom and the CHARMM36 force field

CHARMM (*Chemistry at HARvard Macromolecular Mechanics*) is a MD engine, which uses its own force fields. They are bunched together as CHARMM force fields and are provided for other MD engines as well. The number of the force field corresponds to the version number of CHARMM, in which the force field was used first. The force field CHARMM36 [5, 22] (C36) was published in 2010.

In C36 all atoms are considered (all-atom force field). Parameters were mainly optimized to structural experimental data, such as nuclear magnetic resonance (NMR) or X-ray data, but also QM and semi empirical QM calculations were used (e.g. for dihedral angles of the sidechains of proteins [5] and partial charges of lipids [22]).

All simulations for optimizing have been done with a time step of 2fs. Therefore very detailed dynamics are still included in the simulations.

There are several models for water, which can be used with all-atom force fields. They differ mainly in intermolecular interactions of water atoms. For simulations

in this thesis the TIP3P water model was used. It is a modification of the TIP3P water model [20] used in parametrisation of C36, which allows also Lennard-Jones-interactions on the hydrogens and not only on the oxygen [5].

Coarse graining and the Martini 2.2 force field

The Martini 2.2 [25, 19] force field was introduced in 2013 and is one of the most famous coarse graining force fields. It maps usually four heavy atoms onto a single bead (ring-like structures need a higher resolution), which implicates of course a loss of chemical information, but also an enormous reduction of computational cost. With this approach much larger time- and spatial scales are accessible for MD simulations.

The parametrisation in Martini is mainly based on fitting partitioning free energies of small molecules, such as amino acid side chain analogues, between a range of polar and non-polar solvents to results from all-atom simulations and experimental. For lipids also thermodynamic properties, such as area per lipid, have been considered. The parameters for membrane-protein interactions were optimized by fitting binding energies of small peptides and a membrane to experimental data and results from all-atom simulations data. In the simulations for optimization a time step of 20fs up to 30fs was used [19, 25].

Coarse graining brings a lot of side effects along. First of all, the coarse graining is not unique, therefore important structural properties (e.g. the lipid tail length) are neglected. In proteins the simplifications also lead to problems, because the secondary structure becomes less stable in coarse grained models and has to be constrained by elastic networks (additional bonds between backbone beads) [26, p. 6812].

Furthermore coarse grained beads have a larger size, which leads a.o. to a smoothing of the energy profile. Because smaller local minima in the energy profile, which would slow down the evolving of the system, are smoothed out, coarse graining speeds up the dynamics in the system. The speed up factor is not constant, but can be relatively good approximated by factor four (obtained in most diffusion simulations) [26, p. 6810, 27, p. 7815].

Coarse graining also has an effect on the entropy and the temperature dependency

of the system. In NpT ensembles the Gibbs free energy G is given by

$$G = H - TS \quad (1.10)$$

where S is the entropy and H the enthalpy. Due to a lower number of degrees of freedom in coarse grained systems the configurational entropy is reduced. Because Martini is tuned to free energy calculations, this implicates also a reduction of the enthalpy. Therefore a decomposition of G might not be realistic [26, p. 6811].

In Martini solvent beads represent four water molecules. Unfortunately the Martini water has a freezing temperature of up to 300 K and the freezing process is very sensitive to nucleation. To prevent this antifreeze beads were introduced, which have the same properties as the standard water beads except for a larger σ for the Lennard-Jones potential. By changing around 10% of the standard solvents to these antifreeze beads, the lattice conformation is disturbed and the freezing temperature is increased [27, p. 7815].

The treatment of water as uncharged beads also implies, that the water phase is not polarizable in Martini. This problem is addressed by assuming a uniform relative dielectric constant, but at water phase interfaces the electrostatic interactions are systematically wrong and the interaction strength of polar beads is underestimated. To overcome this a polarizable water (PW) model was introduced to Martini. PW consists of three beads with equal mass. There is one positively charged bead and one negatively charged, both are bound to the central neutral bead. The charged beads acting only electrostatically with other charged beads, but not intramolecular. To control the distribution of the dipole momentum the binding angle between the three beads is constrained. The central bead interacts via the Lennard-Jones potential [34]. However PW comes with a higher computational cost and is not as well tested as the standard water model.

1.3.3 External constraints

Periodic boundary conditions

Because the simulation of an open system is not possible, boundary conditions have to be considered. Closed boundaries often lead to surface interaction artefacts and are therefore in most cases not suitable for MD simulation. That is why usually periodic boundary conditions (PBC) are used.

To use PBC the shape of the simulation box has to have a space filling geometry (e.g. rectangular or rhombic dodecahedron). With periodic boundary conditions images of the simulation box are repeated in every direction. If a particle leaves the simulation box, its periodic image is coming in from the opposite. So the number of particles is kept constant while surface interactions are avoided.

A particle interacts only with the nearest image of another particle, which means that particles near boundaries can interact with periodic images of other particles instead of the real particle. Nevertheless molecules can have long range interaction with their own periodic images, which leads to artefacts and has to be considered when choosing the systems size.

It is generally thought, that PBC only have small or no effects on equilibrium properties or structures of fluids, but known problems are the absence of long wavelength fluctuations (e.g. in phase transitions) or the violation of angular momentum conservation [28, p. 141f].

Thermostats

Integration schemes in MD are designed to conserve the energy of a system. This refers to a microcanonical ensemble, in which the number of particles N , the volume V and the energy E is constant (NVE). However in real biological systems not the energy but rather temperature is kept constant, which is the canonical ensemble (NVT). This can be achieved by thermostats.

In this thesis the Berendsen thermostat, the Parinello-Bussi and the Nosè-Hoover thermostat were used.

The Berendsen thermostat [4] couples the system weakly to a heat bath with temperature T_0 by scaling the velocities of the particles. Weak coupling means, that the energy difference is not conducted in one rescaling, but over a given time scale τ_T , which allows a specification of the coupling strength. As the rescaling of velocities transfer kinetic energy from internal degrees of freedom, such as vibrations, to translational and rotational kinetic energy of the center of mass of the system [17, p. 738], this thermostat is not suitable for production runs. However deviations from T_0 decay exponential making this thermostat usefull for equilibration runs or non equilibrium MD simulations [4, p. 3689].

For the thermostats used in this thesis some characteristics are outlined below.

Berendsen thermostat The Berendsen thermostat [4] couples the system weakly to a heat bath with temperature T_0 by scaling the velocities of the particles. The temperature T , to which the system is set to, is given by

$$\frac{dT}{dt} = \frac{T_0 - T}{\tau_T} \quad (1.11)$$

Therefore an exponential relaxation can be observed. Weak coupling means, that the energy difference is not conducted in one rescaling, but over a given time scale τ_T , which allows a specification of the coupling strength.

It has been shown, that rescaling of velocities transfer kinetic energy from internal degrees of freedom, such as vibrations, to translational and rotational kinetic energy of the center of mass of the system [17, p. 738], which results in a wrong sampling of phase space. Therefore this thermostat is not suitable for production runs. Due to the exponential decay it adjusts quit fast, is stable against large deviations from the desired temperature and prevents oscillations. For these reasons it is often used in equilibration runs or non-equilibrium MD simulations [4, p. 3689].

Parrinello-Bussi thermostat The Parrinello-Bussi thermostat [11] extends the Berendsen thermostat by a stochastic term, which leads to a more accurate sampling of phase space. The advantages of the Berendsen thermostat are still in place [1, p. 31].

Nosé-Hoover thermostat The Nosé-Hoover thermostat [29, 18] extends the Hamiltonian of the system by a friction representing a heat bath. The equations of motion are modified to

$$\frac{d^2 \vec{r}_i}{dt^2} = \frac{\vec{F}_i}{m_i} - \frac{p_\xi}{Q(T_0, \tau_T)} \frac{d\vec{r}_i}{dt} \quad (1.12)$$

Here ξ represents the friction parameter with momentum p_ξ and a mass parameter $Q(T_0, \tau_T)$, which defines the coupling strength. Q depends on the target temperature T_0 and a coupling time scale τ_T . The evolution of ξ is defined as

$$\frac{dp_\xi}{dt} = T - T_0 \quad (1.13)$$

This leads to oscillations between the system and the heat bath with the period τ_T . The relaxation is about five times slower as by an exponential relaxation (see Berendsen or Parrinello–Bussi thermostats) with the same τ_T [1, p. 32f].

Often groups are coupled to independent thermostats. This is helpful, because the heat exchange between f.e. proteins and solvents is often not correct. Therefore proteins would cool down and the solvent would heat up [1, p. 34].

Barostats

In biological systems often not the volume but the pressure is constant. This refers to the isobaric isothermal ensemble (NpT). Below the characteristics of the barostats used in this thesis are outlined.

Berendsen barostat Analogously to the Berendsen thermostat the Berendsen barostat [4] couples the systems pressure P weakly to an external pressure P_0 by rescaling the positions of the particles.

$$\frac{dP}{dt} = \frac{P_0 - P}{\tau_P} \quad (1.14)$$

τ_P is the time scale of coupling. Similar to the Berendsen thermostat this barostat does not sample the NpT ensemble and should therefore only be used in equilibration runs and non-equilibrium MD simulations [1, p. 36].

Parinello-Rahman barostat In the Parinello-Rahman barostat [31, 30] the equations of motion of the particles change to

$$\frac{d^2 \vec{r}_i}{dt^2} = \frac{\vec{F}_i}{m_i} - \underline{M} \frac{d\vec{r}_i}{dt} \quad (1.15)$$

where \underline{M} is a matrix is given by a differential equation depending on the current pressure, the target pressure, the volume and a time scale. The Parinello-Rahman barostat samples the full phase space and can be therefore used in production runs. Large deviations from the desired pressure however lead to oscillations in the box, therefore it should not be used for equilibration runs [1, p. 36].

GROMACS provides the possibility to couple the z-direction independently from the x- and y-direction, which is called semiisotropic pressure coupling. This feature is useful for membrane and pulling simulations, because the dynamics differ a lot between these axes.

1.4 Free energy

To understand state transitions in a physical system the free energy is very handy quantity. It is directly linked to the probability distribution for different states and other quantities can be derived easily.

Below the basic concept of free energy is outlined as well as a practical method to retrieve free energy landscapes from MD simulation, which was used in this thesis.

1.4.1 Partition functions and free energies

The behaviour of a system depends on the interaction with the outside of the system. These interactions are classified in ensembles. The most important are the microcanonical ensemble (N, V and E constant), canonical ensemble (N, V and T constant) and the isothermal-isobaric ensemble (N, p and T constant).

In the microcanonical ensemble the probability, that a system enters a microstate with energy $E' = E(\mathbf{q}, \mathbf{p})$ (\mathbf{q}, \mathbf{p} are the positions and momenta of the particles respectively), is equal for all $E' < |E + dE|$ and 0 else. Therefore the partition function Ω is given as

$$\Omega(N, V, E) = C_0 \int \delta(\mathcal{H}(\mathbf{q}, \mathbf{p}) - E) d\mathbf{q}d\mathbf{p} \quad (1.16)$$

where $\mathcal{H}(\mathbf{q}, \mathbf{p}) = U(\mathbf{q}) + K(\mathbf{p})$ is the Hamiltonian and C_0 a proportional constant, in which the smallest phase space volume and the indistinguishability of particles have to be taken into account [8, p. 16].

In the canonical ensemble however the temperature T is kept constant instead of the energy. Therefore the partition function has to include all possible energies weighted with their probability given by the Boltzmann factor. Below $\frac{1}{k_B T}$ is shortened with β .

$$Q(N, V, T) = \int \exp(-\beta E) \Omega(N, V, E) dE \quad (1.17)$$

$$= C_0 \int \exp(-\beta \mathcal{H}(\mathbf{q}, \mathbf{p})) d\mathbf{q}d\mathbf{p} \quad (1.18)$$

The configurational integral is defined as

$$Z(N, V, T) = \int \exp(-\beta U(\mathbf{q})) d\mathbf{q} \quad (1.19)$$

It is important to see, that \mathcal{H} only depends on the quadrature of \mathbf{p} , so the integral over the momenta can always be solved analytical by turning it into a Gauss integral. This implies, that for two related systems, in which the particle masses are the same, the integral over \mathbf{p} does not change and therefore

$$\frac{Q_2}{Q_1} = \frac{Z_2}{Z_1} \quad (1.20)$$

holds [8, p. 17].

The partition function of the isothermal-isobaric ensemble, in which the pressure is kept constant instead of the volume, can be set up by expanding the Hamiltonian by the work the system is doing on expansion/contraction.

$$\mathcal{H}' = \mathcal{H} + pV \quad (1.21)$$

$$\Xi(N, p, T) = C_1 \int \exp(-\beta pV) Q(N, V, T) dV \quad (1.22)$$

where C_1 is a volume scale. For a further discussion of C_1 please [see 16].

The Helmholtz free energy A refers to a canonical ensemble, which means that it is minimal in equilibrium under constant temperature and volume. The Gibbs free energy G refers to an isothermal-isobaric ensemble [8, p. 19].

$$A = -\beta \ln(Q) \quad (1.23)$$

$$G = -\beta \ln(\Xi) \quad (1.24)$$

Usually the exact value of the free energy is unknown, but the main interest lies in free energy differences between two states of a system. For the Helmholtz free energy A this difference is given as

$$\Delta A = A_2 - A_1 = -\beta \ln(Q_2) + \beta \ln(Q_1) = -\beta \ln\left(\frac{Q_2}{Q_1}\right) = -\beta \ln\left(\frac{Z_2}{Z_1}\right) \quad (1.25)$$

Free energy in MD and Umbrella Sampling

The free energy of a system as a function of a set of parameters $\vec{\xi}$ is given as

$$A(\vec{\xi}) = -\beta \ln \left(\rho(\vec{\xi}) \right) \quad (1.26)$$

$\rho(\vec{\xi})$ can be easily measured during a MD simulation by counting the number of states, in which $\vec{\xi}(\mathbf{q}) = \vec{\xi}$. Therefore it is also referred as histogram. $\vec{\xi}$ are often called reaction coordinates and could be e.g. the distance between two molecules. In reality however $A(\vec{\xi})$ can have big barriers and because the potential energy U is sharply distributed around its mean in NVT simulations, $\vec{\xi}$ can only hardly be sampled during a finite simulation.

One possibility to overcome the sampling problem is umbrella sampling [32]. In this approach the path ξ (for simplicity one dimensional) is split up into distinct windows $[\xi_0, \xi_n]$. To each window i a biasing potential $\hat{U}_i(\xi)$ can be applied to ensure a good sampling of ξ around ξ_i . This changes the potential energy to

$$U_{B,i}(\mathbf{q}) = U(\mathbf{q}) + \hat{U}_i(\xi(\mathbf{q})) \quad (1.27)$$

After a simulation with the biased potential $U_{B,i}$ the unbiased probability distribution $\rho_i(\xi)$ has to be reconstructed from the observed biased one $\rho_{B,i}(\xi)$.

In the following explanation the window index i is dropped and the unbiased system and the biased system are referred as 1 and 2 respectively.

The probability density of finding the biased system in a configuration, in which its potential energy differs from the potential energy of the unbiased system (that is in the same configuration), by $\Delta U = U_2(\mathbf{q}) - U_1(\mathbf{q})$ shall be considered. This implies $U_2(\mathbf{q}) = U_1(\mathbf{q}) + \Delta U$, which is done by an Delta function.

$$\rho_2(\Delta U) = \frac{\int \exp(-\beta U_2(\mathbf{q})) \delta(U_2(\mathbf{q}) - U_1(\mathbf{q}) - \Delta U) d\mathbf{q}}{Z_2} \quad (1.28)$$

By substituting $U_2(\mathbf{q})$ the factor $\exp(-\beta \Delta U)$ can be moved out of the integral. After a multiplication of $\frac{Z_1}{Z_2}$, $\rho_1(\Delta U)$ can be identified. $\rho_1(\Delta U)$ has the same meaning as $\rho_2(\Delta U)$, but refers to the unbiased system. This is the corrected probability density [14, p. 179ff].

$$\rho_2(\Delta U) = \frac{Z_1}{Z_2} \exp(-\beta \Delta U) \rho_1 \quad (1.29)$$

Because $\Delta A = -\beta \ln \left(\frac{Z_2}{Z_1} \right)$, Equation 1.29 can be transformed into

$$\rho_1(\Delta U) = \exp(-\beta(\Delta A - \Delta U)) \rho_2(\Delta U) \quad (1.30)$$

Therefore the measured biased histograms $\tilde{\rho}_{B,i}(\xi)$ can be turned into unbiased histograms $\tilde{\rho}_i(\xi)$ via

$$\tilde{\rho}_i(\xi) = \exp \left(-\beta \left(\Delta A_i - \hat{U}_i(\xi) \right) \right) \tilde{\rho}_{B,i}(\xi) \quad (1.31)$$

Of course ΔA_i is not known, but assuming that $\rho(\xi)$ is a continuous function, the results from the single windows can be combined and afterwards normalized. With this method however one can only have two histograms in the overlapping region. Another problem is, that the sampling in the tails is usually poor and statistical errors, which propagate through all overlapping regions, can become very large [8, 236ff]. Therefore umbrella sampling is usually combined with the Weighted histogram analysis method (WHAM) [13, 23]. This algorithm is able to combine several histograms in one overlapping region and it is designed to keep the statistical errors small. The main idea is to combine probability densities linearly with an additional weighting factor $\omega_i(\xi)$ to the total probability density $\tilde{\rho}$. The weighting factors ω_i are chosen iteratively in a way, that the overall statistical error is minimized.