

# Analyzing and biasing simulations with PLUMED

Giovanni Bussi and Gareth A. Tribello

February 2, 2018

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Collective variables</b>	<b>4</b>
2.1	Ensemble averages . . . . .	5
2.2	Free-energy landscapes . . . . .	6
2.3	Analyzing simulations with PLUMED . . . . .	8
2.4	Distances, angles and torsions . . . . .	10
2.5	Positions and RMSD . . . . .	12
2.6	Gyration radius and gyration tensor . . . . .	15
2.7	Dealing with periodic boundary conditions . . . . .	18
2.8	Going further with collective variables . . . . .	19
<b>3</b>	<b>Biasing collective variables</b>	<b>23</b>
3.1	Reweighting . . . . .	23
3.2	Extracting the free energy . . . . .	24
3.3	Biasing simulations with PLUMED . . . . .	25
3.4	Free Energies . . . . .	27
<b>4</b>	<b>Calculating error bars</b>	<b>28</b>
4.1	Block analysis with PLUMED . . . . .	32
<b>5</b>	<b>Multiple replica techniques</b>	<b>33</b>
5.1	The weighted histogram method . . . . .	35
5.2	WHAM analysis with PLUMED . . . . .	36
<b>6</b>	<b>Perspectives</b>	<b>42</b>
<b>7</b>	<b>Notes</b>	<b>43</b>

# 1 Introduction

The chapters in sections I to IV will have given you some sense of the broad range of methods and techniques that have been used to simulate biomolecular processes. The aim of this chapter is not to introduce more techniques but rather to focus on how these techniques can be employed in practice. We will do so by explaining how a particular piece of software, PLUMED [1, 2], can be used to run and analyze many of the types of simulation that are discussed in section II. Note 1 discusses the various different versions of PLUMED. The most important thing to know at the outset, however is that PLUMED is not a molecular dynamics (MD) or Monte Carlo code. It is instead designed to complement MD codes such as GROMACS [3], LAMMPS [4], DL-POLY [5], CP2K [6], AMBER [7], and OpenMM [8]. PLUMED does this in two ways:

1. It can be used to post-process the molecular dynamics trajectories that are generated by these code.
2. It can serve as a plugin to these MD codes and thus allow the user to add the additional biasing forces that are required for the enhanced sampling methods, such as umbrella sampling or metadynamics, that are described in Section II.

The manner in which PLUMED is plugged into an MD code is illustrated in figure 1. As you can see PLUMED is called during initialization and its input file is read in at that time. It is then called during each run through of the main loop of the MD code just after the forces that describe the interactions between the atoms are calculated. Calling PLUMED at these points allows the plugin to do any analysis that is required and also allows any forces due to bias potentials that are calculated by PLUMED to be returned from PLUMED to the MD code so that they can be incorporated when the equations of motion are integrated. PLUMED is not the only piece of software that interacts with other MD codes in this manner. Two other notable examples are the COLVARS package [9], which is also reviewed in this book, and the recently published program SSAGES [10].

When the PLUMED package is used to post process trajectories a program that is part of PLUMED and that is called **driver** is employed. When PLUMED **driver** is used the trajectory is not generated by integrating the equations of motion. The trajectory is instead read from disk so the forces calculated within PLUMED thus play no role in the systems' dynamics. PLUMED **driver** is, nevertheless, useful as the PLUMED code contains numerous analysis tools that can also be used to post process trajectories.

PLUMED is written using C++ and the object oriented paradigm is heavily used in the design of the code. The consequence of this is that the code contains a set of core modules that look after the communication with the MD codes and various other mission critical features. This code is maintained by a small cadre of core developers. Additional functionalities can be built on this core and contributed by any user in the community. In fact, developers who wish to contribute to the project in this way often only need to contribute one single file

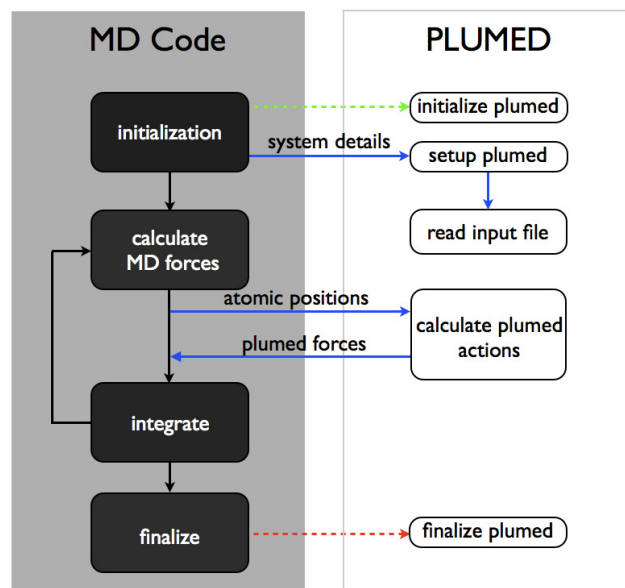


Figure 1: Schematic representation of the interface between PLUMED and an MD engine. The green arrow indicates that the function called should create the PLUMED object and the red arrow indicates that the function called should delete the PLUMED object

that contains class and method definitions, the code itself and the sections of the manual that describe the new functionality and files for a regression test can be used to ensure that the new functionality continues to work if changes are made elsewhere. The fact that it is relatively easy to extend PLUMED ensures that the code, the associated website and the various user meetings provide a forum that developers can use to share new techniques and methods with the scientific community. In fact, in PLUMED v2.4, a modular framework was introduced that allows developers to contribute groups of functionalities that are logically connected. So far three groups of non-core developers have used this model to contribute modules to the code base, but we expect this number to increase in the future.

The size and scope of the PLUMED code ensures that we cannot describe everything that it can do in this single chapter. For those who are interested we would recommend reading the original article [2], the code’s online manual, the many tutorials included in the manual, and some of the considerable number of papers that describe simulations done or analyzed with PLUMED. What we will do in the following is quickly summarize the theory behind some of the methods that are implemented in PLUMED. We will also provide relevant examples of PLUMED input files that can be used for these types of calculations. Our focus in these sections is on describing how free energies are estimated from these types of calculations, how results from multiple replicas can be combined, and how suitable error bars on these estimates are determined.

## 2 Collective variables

Complex biochemical reactions or conformational changes are often interpreted in terms of free-energy surfaces that are computed as a function of a small number of collective variables (CVs). These potentials of mean force provide a coarse grained representation of the energy landscape for the system, which can in turn provide useful insights into the behavior of these biochemical systems. To see why consider the following examples:

- We know that nucleosides have multiple isomers and that isomerization involves a rotation around the glycosidic bond. A free energy surface for such a system as a function of the torsional angle  $\chi$  around this bond is thus useful as it provides us with information on the relative free energies of the two isomers and an indication of the height of the barrier to isomerization (see, e.g., Ref. [11]).
- We know that the secondary and tertiary structures of many proteins determine their function and that it is important to understand the mechanism by which proteins fold. If we run molecular dynamics simulations we can understand something about this folding process by extracting the free energy surface along a coordinate that counts the number of native contacts that are present (see, e.g., [12, 13]).

- We know that cells are constantly exchanging water and ions across their membrane. To extract information on the mechanisms of these processes we can run molecular dynamics simulations and extract a free energy surface that describes how the free energy of the ion changes as the ion moves through an ion channel (see, e.g., [14]).
- We know that the behavior of biomolecules changes when their protonation state changes. To extract information on the likely protonation state of a biomolecule we might, therefore, calculate how the free energy of the molecule changes as the distance between the hydrogen atoms and the various protonation sites on the molecule changes (see, e.g., [15]).
- We know that the solute molecules in a solution must all aggregate in one place in order for a crystal to form. To investigate the ease with which a crystal forms from a particular solution and the earliest stages of this nucleation process we might, therefore, calculate the free energy as a function of the number of molecules that are in the solid phase (see, e.g., [16, 17]).

In all these examples the Hamiltonian depends on numerous degrees of freedom, but we are only interested in the behavior of a small number of these degrees of freedom (e.g., a torsional angle or a distance between two atoms). A CV,  $s$ , is thus simply an arbitrary function of the atomic coordinates,  $q$ , although on some rare occasions it might be a function of velocities as well [18]. As we will see in what follows, in many cases the function  $s$  is very simple and thus easy to calculate from  $q$ . In other cases, however, the function is considerably more complicated so for these cases the scripting language of the PLUMED input file is invaluable.

## 2.1 Ensemble averages

Section IV discussed a variety of different ways in which MD trajectories can be analyzed. The simplest way to analyze the values a CV takes during a trajectory is to compute an ensemble average. We are able to compute such averages from MD simulations because, if we are running at temperature  $T$  on a system containing  $N$  atoms that interact through a Hamiltonian,  $H(\mathbf{q}, \mathbf{p})$ , the probability,  $P(\mathbf{q}, \mathbf{p})$ , that a microstate in which the atoms have positions  $\mathbf{q}$  and momenta  $\mathbf{p}$  will be sampled at a particular instant in time is given by:

$$P(\mathbf{q}, \mathbf{p}) = \frac{\exp\left(-\frac{H(\mathbf{q}, \mathbf{p})}{k_B T}\right)}{\int d\mathbf{q}' d\mathbf{p}' \exp\left(-\frac{H(\mathbf{q}', \mathbf{p}')}{k_B T}\right)} \quad (1)$$

where  $k_B$  is Boltzmann's constant and where the  $6N$ -dimensional integral in the denominator runs over all the possible values of position and momentum that each of the atoms in the system might have. Consequently, an observable  $A(\mathbf{q})$

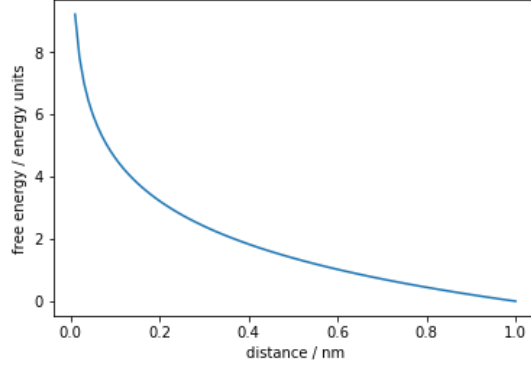


Figure 2: Figure showing the free energy surface as a function of the distance between two atoms that do not interact. As you can see the free energy decreases as the two atoms move apart but this is only because the phase space volume that is accessible to two atoms that are exactly a distance  $r$  apart is proportional to  $r^2$ .

will have an ensemble average that is equal to:

$$\langle A \rangle = \int d\mathbf{q} d\mathbf{p} A(\mathbf{q}) P(\mathbf{q}, \mathbf{p}) \quad (2)$$

By the law of large numbers, however, we know that if we add together the value  $A(\mathbf{q})$  took in each of the frames in our molecular dynamics trajectory and if we divide this sum by the number of frames in our trajectory we will obtain an estimate for  $\langle A \rangle$ .

## 2.2 Free-energy landscapes

A slightly more advanced way to analyze CVs is to compute a marginal probability distribution. This function is defined as:

$$P(s) \propto \int d\mathbf{q} d\mathbf{p} P(\mathbf{q}, \mathbf{p}) \delta(s(\mathbf{q}, \mathbf{p}) - s) \quad (3)$$

The free energy surface,  $F(s)$ , is then nothing more than the negative logarithm of this probability density function expressed in units of energy:

$$F(s) = -k_B T \ln \left[ \int d\mathbf{q} d\mathbf{p} P(\mathbf{q}, \mathbf{p}) \delta(s(\mathbf{q}, \mathbf{p}) - s) \right] + C \quad (4)$$

where  $C$  is an arbitrary constant.

Recognizing this connection between the value of the thermodynamic potential and the probability of having a particular value for a collective variable is fundamental in terms of understanding what a free energy landscape means.

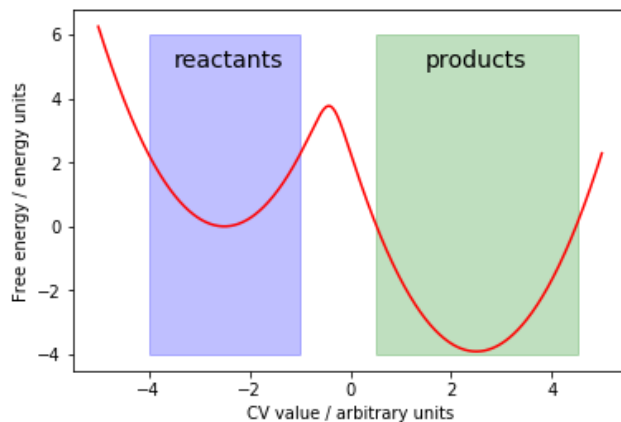


Figure 3: Figure showing a typical one-dimensional free energy surface that we might extract from an MD simulation. It is clear from this diagram that there are two metastable basins in this landscape and that the system must cross a substantial barrier in order to get from one to the other. The two shaded regions in this figure indicate the two sets of limits that we would integrate over when evaluating the free energy difference between these two metastable basins using equation 5.

As a case in point consider the free energy surface shown in figure 2. This figure shows how the free energy changes as the distance between two particles is increased. One might be tempted, based on the shape of the curve, to assume that the two particles repel each other. In actual fact, however, there is no interaction between the particles. The free energy decreases because all possible distance vectors are equally likely. In three dimensions the probability of observing a particular distance,  $r$ , is thus proportional to  $r^2$ . The free energy for a pair of non-interacting particles as a function of the distance between them is thus  $F(r) = -2k_B T \ln r + C$ . Hence, when free energy landscapes are used to provide information on the strength of the interaction between two particles the surface obtained is usually corrected by adding  $2k_B T \ln r$  so that the free energy surface for a pair of non-interacting particles appears flat.

Let us now suppose that we have extracted a free energy surface that looks something like the one shown in figure 3. This free energy surface contains two metastable basins, which we will, for the time being, assume correspond to the reactant and product states for a chemical reaction that our simulated system can undergo. If we want to calculate the free energy change associated with this reaction we would use the expression below:

$$F_{\text{products}} - F_{\text{reactants}} = -k_B T \log \left[ \frac{\int_{\text{products}} e^{-\frac{F(s)}{k_B T}} ds}{\int_{\text{reactants}} e^{-\frac{F(s)}{k_B T}} ds} \right] \quad (5)$$

Here the integrals in the numerator and the denominator of the quotient on the right-hand side run over the regions highlighted in figure 3. Notice, furthermore, that the free energy difference between the reactant and product states is calculated in this way because the value of the CV will fluctuate when it is in either state. Calculating the reaction free energy using the formula above incorporates the effect of these fluctuations whereas simply calculating the difference between the values of the free energy at the bottom of two minimums does not.

A further interesting thing that is worth noting about free energy surfaces is that, if one has the free energy,  $F(s)$ , as a function of some CV,  $s$ , and if one wishes to determine the free energy as a function of some second CV,  $\zeta$ , as long as  $\zeta$  is a one-to-one function of  $s$  with an inverse function that is differentiable one can write an analytic expression for  $F(\zeta)$  in terms of  $F(s)$ . The reason that this is possible is that if  $\zeta(s)$  is one-to-one then we know that:

$$P(\zeta)d\zeta = P(s)ds = P(s) \left| \frac{ds}{d\zeta} \right| d\zeta \quad (6)$$

Once we recall the connection between free energy and probability we thus arrive at:

$$F(\zeta) = F(s) - k_B T \ln \left| \frac{ds}{d\zeta} \right| \quad (7)$$

$F(\zeta)$  may look rather different to  $F(s)$  but these two free energy surfaces will still convey the same information about depths of the basins and the heights of the barriers. In general, however, one should be particularly careful when discussing the barriers found in free energy landscapes, as the height of the barrier will depend on the particular combination of CVs that are used to display the free energy landscape. As illustrated in figure 4, if the critical degree of freedom that distinguishes between two metastable states is in a direction that is orthogonal to the CVs then these two states will both contribute their probability density to a single basin in the final free energy landscape that is extracted. If you want to extract information on barriers you should thus probably also extract the average passage times between states as this better characterizes the behavior of the chemical system. A free energy landscape alone can only tell you about the relative stabilities of the various metastable states that can be distinguished by the CVs that were used in its construction. It cannot really provide information on the rates with which the system will pass from state to state.

### 2.3 Analyzing simulations with PLUMED

In this section we will discuss how to compute various CVs using PLUMED. It is important to emphasize at the outset that the PLUMED input files we provide can be used when running a simulation using some other MD code combined with PLUMED or when analyzing a simulation trajectory using PLUMED **driver**. In the first case the command line input will depend on which precise MD code one is using. For instance, simulations done with GROMACS are typically launched using a command such as:



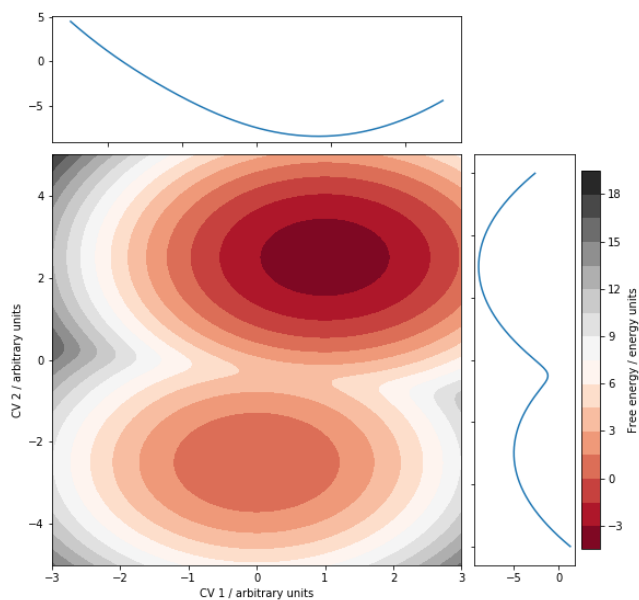


Figure 4: The central panel of this figure shows a two-dimensional free energy surface with two metastable states that is a function of two CVs called CV1 and CV2. The panels above and to the right of this figure show the free energy surface as a function of CV1 and CV2 respectively. It is clear from these two figures that while CV2 is able to distinguish the two metastable basins CV1 is not.

```
> gmx mdrun -plumed plumed.dat
```

where `plumed.dat` is the name of the PLUMED input file. In the second case, however, as only PLUMED is being used, the command is nearly always the same and will be something akin to:

```
> plumed driver --plumed plumed.dat --igro traj.gro
```

where `traj.gro` is a trajectory file in GROMACS format and `plumed.dat` is once again a PLUMED input file.

For the final end user a PLUMED input file looks like it is written in a rudimentary scripting language. Each line in the input file tells the code to do something, which may be as simple as calculating a position of a center of mass or as complex as calculating and accumulating a metadynamics bias. As these commands can be used in a wide range of different contexts and orders the user has the flexibility to do all range of the analyses described in the previous sections. Furthermore, if they choose to incorporate some new functionality they can quickly start to use it in tandem with the methods developed by other members of the PLUMED community.

## 2.4 Distances, angles and torsions

Chemical reactions are one class of phenomena that we can investigate using MD simulations. In many chemical reactions two atoms or groups of atoms approach each other so that a chemical bond can form between them. The ideal CV to use to describe this phenomenon is thus the distance between the relevant atoms. We can calculate and print the distance between a pair of atoms using the PLUMED input file below:

```
d1: DISTANCE ATOMS=1,2
PRINT ARG=d1 FILE=colvar STRIDE=10
```

This input instructs PLUMED to calculate the distance between the positions of the first and second atoms in the MD code's input file and to print this quantity to a file called `colvar` every 10 steps. We can take the output from this calculation and plot a graph showing the value of the distance as a function of time. The resulting curve would look something like the red curve in figure 5.

We can also use PLUMED to calculate the average value of the distance between two atoms. The following input calculates the distance every 10 steps and then calculates a sample mean over the whole trajectory.

```
d1: DISTANCE ATOMS=1,2
a1: AVERAGE ARG=d1
PRINT ARG=a1 FILE=colvar STRIDE=100
```

The output generated by this calculation is plotted in blue in figure 5. There are multiple points in the output file as the above input instructs PLUMED to calculate an ensemble average from the first 100 trajectory frames, the first 200

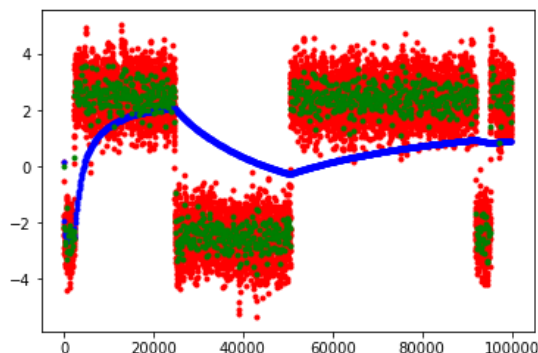


Figure 5: Figure showing some things PLUMED can do with the CVs it calculates. As discussed in the main text to calculate the curves shown in this figure we calculated the distance between atoms 1 and 2 on every 10th MD step. The red points are thus the values of these distances. The blue points indicate the mean value of this distance calculated for a range of differently sized samples. Lastly, the green points are running averages from the first 100 frames, the second 100 frames and so on.

trajectory frames, the first 300 trajectory frames and so on. It is worth noting that one can also use PLUMED to calculate the running averages from the first 100 frames, the second 100 frames and so on - points that are shown using green dots in figure 5 - by using an input file like the one below:

```
d1: DISTANCE ATOMS=1,2
a1: AVERAGE ARG=d1 CLEAR=100
PRINT ARG=a1 FILE=colvar STRIDE=100
```

To calculate and print an angle between two bonds using PLUMED one might use an input file like the one below:

```
ang: ANGLE ATOMS=1,2,3
PRINT ARG=ang FILE=colvar STRIDE=5
```

As shown in figure 6(a) this input calculates the angle between the vector connecting atom 2 to atom 1 and the vector connecting atom 2 to atom 3. The assumption here is that there are chemical bonds connecting atom 2 to atoms 1 and 3. This input is thus measuring the angle between these bonds. As illustrated in figure 6(b) we can, however, also use PLUMED to calculate the angle between any pair of distance vectors. For example the input file below would calculate the angle between the vector that connects atom 2 to atom 1 and the vector that connects atom 3 to atom 4.

```
ang2: ANGLE ATOMS=1,2,3,4
PRINT ARG=ang2 FILE=colvar STRIDE=2
```

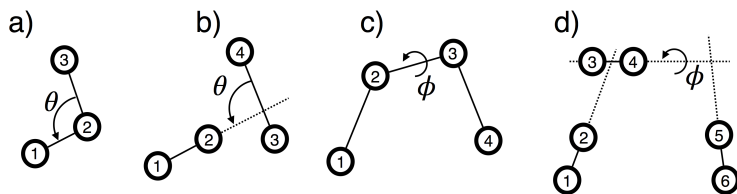


Figure 6: Schematic representations of the angles and torsions that PLUMED can compute. a) Angle defined using three atoms; b) Angle defined using four atoms; c) Torsion defined using four atoms; d) Torsion defined using six atoms.

PLUMED can also be used to calculate and print the torsional angles illustrated in figure 6(c) that involve sets of four atoms as shown below:

```
tor: TORSION ATOMS=1,2,3,4
PRINT ARG=tor FILE=colvar STRIDE=1
```

Much like the command to calculate the angle where three atoms are specified the assumption here is that if the atom 1, 2, 3 and 4 are specified then there are chemical bonds between atoms 1 and 2, atoms 2 and 3 and atoms 3 and 4. In general, however, a torsional angle measures the angle between two planes, which have at least one vector in common as illustrated in figure 6(d). As shown below, there is thus an alternate, more general, way through which we can define a torsional angle:

```
tor2: TORSION VECTOR1=1,2 AXIS=3,4 VECTOR2=5,6
PRINT ARG=tor2 FILE=colvar STRIDE=20
```

This input instructs PLUMED to calculate the angle between the plane containing the vector connecting atoms 1 and 2 and the vector connecting atoms 3 and 4 and the plane containing this second vector and the vector connecting atoms 5 and 6. Notice that there is one additional syntax for selecting the atoms that should be used to calculate a torsion that is discussed in Note 2.

## 2.5 Positions and RMSD

One of the processes that we stated we might want to study was the diffusion of ions across cell membranes. If the membrane is parallel to the  $xy$  plane, one might be tempted to conclude that that correct CV to use in this case is the  $z$  component of the position of the ion. Assuming that the ion is atom number 1001, this can be done with the following input file

```
pos: POSITION ATOM=1001
PRINT ARG=pos.z FILE=colvar STRIDE=10
```

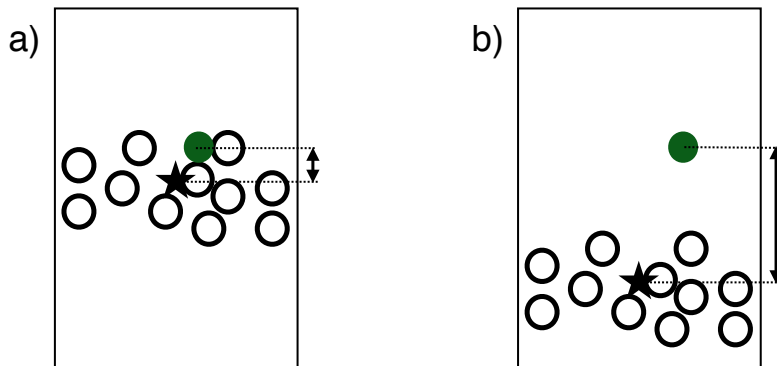


Figure 7: Figure illustrating the vertical position of an ion with respect to a membrane. The ion is shown in green, whereas the atoms of the membrane are depicted as empty circles. The center of the membrane is represented using a star. Panels a) and b) report two different structures. Notice that the vertical coordinate of the ion is the same, but that its position relative to the membrane is different. The correct way to describe the position of the ion with respect to the membrane is thus to use the vertical component of the distance between the ion and the center of the membrane.

The variable `POSITION` that is used here has multiple components. As we want the `z` component specifically, we thus use `pos.z` in the `PRINT` command. Using this CV is a bad idea as the potential energy functions that we use within MD simulations are almost always translationally invariant. As illustrated in figure 7 for our membrane example a particular value for the `z` component of our ion's coordinate might be inside the membrane at one particular time. During the simulation, however, the membrane can move and as such this same value for the ion's `z` coordinate might be outside the membrane at some later time. It is thus usually much better to use the `z` component of the distance between the coordinates of the ion and the center of the membrane and to thus track the `z` position of the ion relative to the `z` position of the membrane. An example input that can be used to calculate and print this quantity is given below:

```
c1: COM ATOMS=1-1000
dist: DISTANCE ATOMS=c1,1001 COMPONENTS
PRINT ARG=dist.z FILE=colvar STRIDE=10
```

Here we assumed that the first 1000 atoms together form the membrane.

It is worth noting here how the `COM` command is used to keep track of the position of the center of mass of the large number of atoms that make up the membrane and how the position of this center of mass is referred to in the `DISTANCE` command using the label, `c1`, of the `COM` command. It is also important to understand how `PLUMED` deals with periodic boundary

conditions (see Section 2.7) and to remember that, for the position of the center to be computed correctly, the vertical span of the membrane must be less than half of the box height. Finally, one should be aware that using a large number of atoms when calculating a CV may well slow down the calculation. It may thus be worth using only a subset of the atoms in the membrane when calculating the position of the center by, for instance, replacing the first line of the input file above with:

```
c1: COM ATOMS=1-1000:10
```

The extra `:10` here tells PLUMED to only use every 10th atom in the range specified. This is a crude solution, however, and it is probably always better in practice to select specific atoms using your understanding of the chemistry of the system.

It can be useful to consider the atomic positions directly when considering problems such as protein folding. Let’s suppose that you know the precise arrangement that the atoms have in the folded structure and you want to monitor the progression of folding during a trajectory. An obvious CV to measure would be the degree of similarity between the instantaneous coordinates of the protein and this special, folded configuration. A method that is commonly used to calculate this degree of similarity involves computing the root-mean-square deviation (RMSD) between the instantaneous coordinates and the coordinates in the reference structure using:

$$s = \sqrt{\frac{1}{N} \sum_{i=1}^N \left( \mathbf{q}_i - \mathbf{q}_i^{(\text{ref})} \right)^2} \quad (8)$$

In this expression  $N$  is the number of atoms,  $\mathbf{q}_i$  is used to denote the instantaneous coordinates of atom  $i$  and  $\mathbf{q}_i^{(\text{ref})}$  is used to denote the coordinates of atom  $i$  in the reference structure. It is important to note that if you were to use the formula above you would have the same problem as if you used the position of an atom as a CV. Consequently, the RMSD formula above is usually computed from a reference frame that is found by performing translation and rotation operations on the original reference structure that minimize the value of the RMSD [19]. An input file that can be used to calculate and print this quantity using PLUMED is shown below:

```
rmsd: RMSD REFERENCE=ref.pdb TYPE=OPTIMAL
PRINT ARG=rmsd FILE=colvar STRIDE=5
```

This input computes the root-mean-square deviation between the instantaneous positions of the atoms that are listed in the input `ref.pdb` file and the positions of those atoms in the `pdb` file. In this case any translation of the center of mass and rotation of the reference frame is removed before calculating the displacements that enter equation 8. As discussed in Note 3, however, PLUMED provides a range of options that allow you to perform these RMSD calculations in various different ways.

As we will see, PLUMED is often used to apply a simulation bias that is a function of the CVs it computes. The forces due to this bias are then propagated onto the atoms that were used to calculate the CV. When these CVs are calculated using the RMSD procedure outlined above the forces required to restore invariance with respect to translations and/or rotations must be applied to all the atoms employed in the alignment procedure, which is computationally expensive. In fact even if you don't have any forces to propagate just performing the alignment operation with a large number of atoms is expensive. One would, therefore, typically never use the positions of all the atoms when performing RMSD calculations.

An alternative to RMSD is the so-called DRMSD:

$$s = \sqrt{\frac{1}{M} \sum_{ij} \left( d_{ij} - d_{ij}^{(\text{ref})} \right)^2} \quad (9)$$

In this expression the sum runs over  $M$  pairs of atoms and  $d_{ij}$  is used to denote the instantaneous distance between these pairs while  $d_{ij}^{(\text{ref})}$  is used to denote the distance between the corresponding pair of atoms in the reference structure. An input file that calculates and prints this quantity using PLUMED is given below. Notice how the two cutoff keywords can be used to specify the range of values the distances should take in the reference structure in order to be considered as part of the sum in equation 9.

```
drmsd: DRMSD REFERENCE=ref.pdb LOWER_CUTOFF=0.1 UPPER_CUTOFF=0.8
PRINT ARG=drmsd FILE=colvar STRIDE=5
```

## 2.6 Gyration radius and gyration tensor

The CVs that have been discussed thus far are all based on a very detailed view of the positions of particular atoms. CVs that give a more coarse-grained view of the instantaneous shape of a molecule can also be used, however. One particularly popular CV of this type is the so-called radius of gyration, which can be calculated using:

$$s = \sqrt{\frac{1}{\sum_i w_i} \sum_i w_i \left( \mathbf{q}^{(i)} - \mathbf{q}^{(c)} \right)^2} \quad (10)$$

where  $\mathbf{q}^{(c)} = \frac{\sum_i w_i \mathbf{q}^{(i)}}{\sum_i w_i}$  and  $\mathbf{q}^{(i)}$  is the position of the  $i$ th atom. The  $w_i$  are a set of weights that are ascribed to each of the atoms in the system. These weights might be, for instance, the masses of the atoms. To calculate and print this quantity using PLUMED you would use an input like the one below:

```
gyr: GYRATION TYPE=RADIUS ATOMS=10-20
PRINT ARG=gyr STRIDE=1 FILE=colvar
```

This calculates the radius of gyration using the positions of the 10th to the 20th atom in the MD code's input file and sets the weights of all these atoms equal to 1.

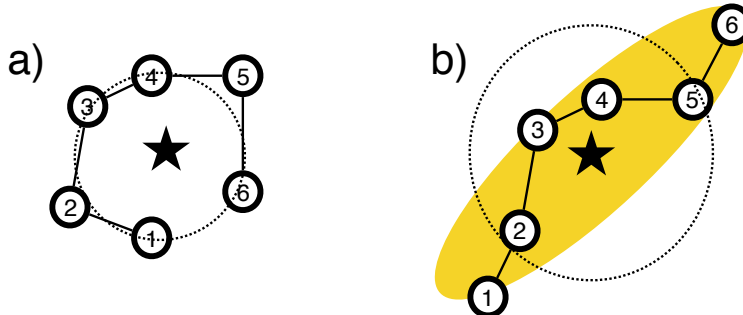


Figure 8: Figure illustrating the components of the radius of gyration of a polymer. The center of the polymer is represented as a star. A circle with a radius corresponding to the radius of gyration of the polymer and centered on the center of the polymer is shown using a dashed line. a) For a globular polymer, the circle correctly represents the shape of the polymer. b) For an elongated polymer, however, the circle does not represent the shape of the polymer correctly. An ellipsoid with axes that have lengths corresponding to the square roots of the eigenvalues of the gyration tensor is more representative and is shown in yellow.

When the gyration radius is calculated using equation 10 the quantity output provides a measure of the average radius of the molecule. No information on the shape of the molecule is provided, however, and so this average radius may be misleading. For example, and as shown in figure 8, the radius of gyration for an extended polymer, which has a shape that is very anisotropic, would not be representative of the extent of the molecule in either its extended or compact directions. For this reason, some researchers have chosen to use the eigenvalues of the gyration tensor instead. The elements of the  $3 \times 3$  gyration tensor are computed using:

$$s_{jk} = \frac{1}{\sum_i w_i} \sum_i w_i \left( q_j^{(i)} - q_j^{(c)} \right) \left( q_k^{(i)} - q_k^{(c)} \right) \quad (11)$$

where  $q_j^{(i)}$  is used to denote the  $j$ th component of the position of atom  $i$  and  $q_j^{(c)}$  is used to denote the  $j$ th component of  $\mathbf{q}^{(c)}$ . Figure 8 shows how the square roots of the eigenvalues of this matrix give a sense of the shape of the molecule. By changing the word after the TYPE keyword in the PLUMED input above one can access these eigenvalues directly or various functions of these eigenvalues that can be used to give one a sense of the sphericity of the molecule or the cylindricity [20].



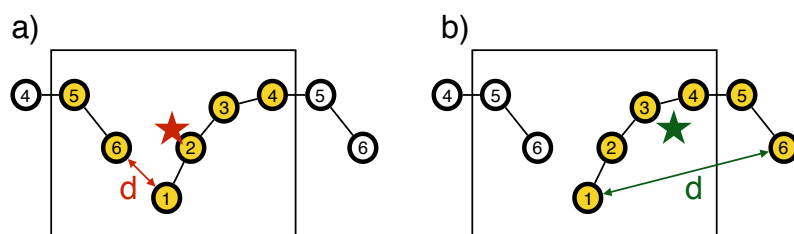


Figure 9: An illustration showing how molecules can be split by the periodic boundary conditions and how this can cause a problem when computing collective variables. Panels a) and b) represent the same set of six atoms. In both of these figures these atoms sit in the periodic box indicated using the rectangle so periodic images of atoms 4, 5, and 6 are included. The atoms used to calculate the collective variables are highlighted. In panel a), a broken molecule is used so the position of center of mass (indicated using a star) and the end-to-end distance (indicated using an arrow) are computed incorrectly. In panel b), however, the molecule has been correctly reconstructed across the periodic boundaries using `WHOLEMOLECULES` and the center of mass has thus been computed correctly. Notice that the correct value for the end-to-end distance is only obtained from panel (b) if the periodic boundary conditions are *ignored* when computing the distance (keyword `NOPBC`). If the PBC are taken into account the incorrect image of atom 6 will be used and the incorrect result illustrated in panel (a) will be obtained once more.

## 2.7 Dealing with periodic boundary conditions

Figure 9 illustrates a technical problem that can appear when PLUMED is used to calculate some CVs. When the underlying MD code applies the periodic boundary conditions molecules can end up split across either side of the periodic box. Consequentially, atoms can appear to be much farther apart than they are in actuality and, as shown in the figure, the position of the center of mass of the molecule can be calculated wrongly. In this case, however, and in some of the others discussed thus far PLUMED resolves this problem automatically by adjusting the positions so that the set of molecules that are used to calculate the position of a center of mass or a CV form one single unbroken molecule. This is still a problem that any user must be aware of, however, as there are some cases that PLUMED cannot fix automatically. For example suppose that you wanted to calculate the end to end distance for the molecule illustrated in 9. In order to do this correctly one must reconstruct the whole molecule before calculating the distance between the two terminal atoms. To resolve this problem PLUMED provides a command, WHOLEMOLECULES, that allows one to adjust the way the positions are stored and to thus specify the molecules that must be reconstructed. A sample input that calculates this end to end distance and that uses a WHOLEMOLECULES command is provided below

```
WHOLEMOLECULES ENTITY0=1-6
d1: DISTANCE ATOMS=1,6 NOPBC
```

The WHOLEMOLECULES command here ensures that the bond between each pair of adjacent atoms specified to the ENTITY keyword is not broken by the periodic boundaries. The distance is thus computed from the positions of the atoms shown in the right panel of figure 9. There is thus no need to apply periodic boundary conditions. In fact if, when the molecule is extended, it has a length that is longer than half the box length it is wrong to apply periodic boundary conditions as the “end-to-end distance” computed this way would no longer be representative of the distance along the chain.

A more complicated example which has been taken from Ref [21] and which also requires the WHOLEMOLECULES command to be used is shown in figure 10. In this case the reconstruction has been done using the following input file

```
MOLINFO STRUCTURE=ref.pdb

rna: GROUP ATOMS=1-258
mg:  GROUP ATOMS=6580
wat: GROUP ATOMS=259-6579

# Make the RNA duplex whole.
WHOLEMOLECULES ENTITY0=rna

# Align the RNA duplex to a reference structure
FIT_TO_TEMPLATE REFERENCE=ref-rna.pdb TYPE=OPTIMAL
```

```

# Notice that before using FIT_TO_TEMPLATE
# we used WHOLEMOLECULES to make the RNA whole
# This is necessary otherwise you would be aligning
# to a broken molecule!

# compute the center of the RNA molecule
center: CENTER ATOMS=rna

# Wrap atoms correctly around the center of the RNA
WRAPAROUND ATOMS=mg AROUND=center
WRAPAROUND ATOMS=wat AROUND=center GROUPBY=3
# Dump the resulting trajectory
DUMPATOMS ATOMS=rna,wat,mg FILE=rna-wrap.gro

```

Notice here how the action `FIT_TO_TEMPLATE` has been used to align the RNA molecules to a template structure that is at the center of the box and how the action `WRAPAROUND` has been used to reposition the water molecules around the aligned RNA molecule. In addition, notice that by aligning the molecule to a template we have made the positions of these atoms roto-translationally invariant. We can thus use the positions of the aligned atoms as CVs using the `POSITION` command that was introduced earlier. Furthermore, when we do so the `FIT_TO_TEMPLATE` command will ensure that the process of aligning the molecule is considered correctly when propagating any forces to the underlying positions.

## 2.8 Going further with collective variables

Hopefully the example input files in the previous sections have given you a sense of how PLUMED input files work. In essence, all the PLUMED commands for calculating CVs that we have introduced thus far calculate a scalar valued quantity. The scalar valued quantity that we calculate in these commands can then be referred to later in the input file using the label of the command that calculates it. So for instance in the following PLUMED input file:

```

dist: DISTANCE ATOMS=1,2
PRINT ARG=dist FILE=colvar STRIDE=1

```

the `PRINT` command instructs PLUMED to print the quantity called `dist` which is calculated by the first `DISTANCE` command to a file called `colvar` during every MD step. The fact that we can pass scalar valued quantities in this way is enormously useful as it means we can script new CVs directly from the input file. For example suppose that we want to calculate the number of native contacts in a protein. This CV is often computed using the continuous switching function shown below or by using some other function that displays a similar behavior [12]:

$$s = \sum_{ij \in N_C} \frac{1}{1 + (r/r_0)^6} \quad (12)$$

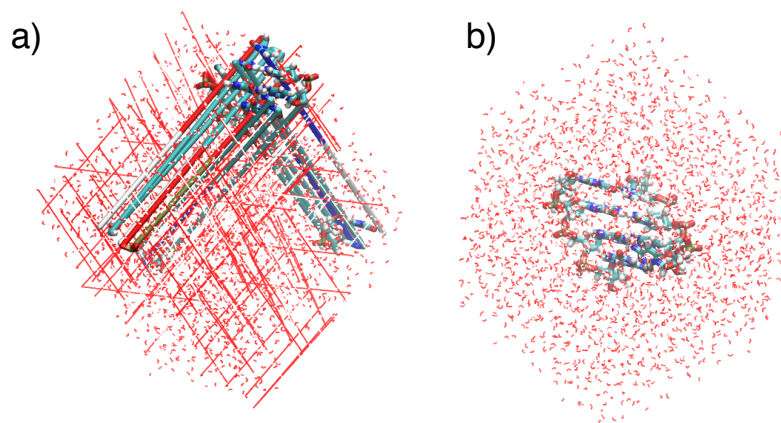


Figure 10: Figures showing how PLUMED can be used to reconstruct a solvated RNA molecule that has been split by the periodic boundary conditions. The RNA duplex is represented using licorice and the water molecules are represented using lines. In panel (a) the molecules that cross the periodic boundaries are broken. In panel (b), however, the RNA molecule has been reconstructed and centered, and the periodic images of the water molecules that are closest to the center of the RNA molecule have been selected. Notice that selecting the water molecules in this way ensures that the water molecules are made whole.

Here the sum runs over the list of pairs of atoms that are in contact when the protein is folded and  $r_0$  is a parameter. If we had a total of four native contacts in the protein and if  $r_0$  were set equal to 6 Å we could compute and print this quantity using the following PLUMED input file.

```
d1: DISTANCE ATOMS=1,2
d2: DISTANCE ATOMS=5,6
d3: DISTANCE ATOMS=9,10
d4: DISTANCE ATOMS=15,16
# the dots here indicate that the command
# will be continued on the following line.
# The braces, meanwhile, allow us to use spaces
# within the argument of the FUNC keyword
contacts: CUSTOM ...
  ARG=d1,d2,d3,d4 VAR=a,b,c,d
  FUNC={
    1/(1+(a/0.6))^6
    +1/(1+(b/0.6))^6
    +1/(1+(c/0.6))^6
    +1/(1+(d/0.6))^6
  }
  PERIODIC=NO
...
PRINT ARG=contacts FILE=colvar STRIDE=10
```

Notice here how we have used our familiar friend the `DISTANCE` command to calculate each of the distances required and have then used the command `CUSTOM` to calculate the non-linear combination of quantities that the CV requires. The fact is that most of the more complicated CVs that have been used to analyze molecular dynamics trajectories are simply linear or non-linear combinations of the quantities that have been introduced thus far. One can thus compute many complicated CVs by just using the commands that were introduced in the previous sections together with the `CUSTOM` command. It is worth understanding how to use this approach, but we would, in practice, not recommend you write your PLUMED input files in this way. Instead, we would recommend that you use the numerous shortcuts PLUMED provides for accessing these non-linear combinations. For example a shorter input that does the same calculation as the input file above is as follows:

```
dists: DISTANCES ...
  ATOMS1=1,2 ATOMS2=5,6 ATOMS3=9,10 ATOMS4=15,16
  LESS_THAN={RATIONAL R_0=0.6}
...
PRINT ARG=dists.less_than FILE=colvar STRIDE=10
```

This input takes advantage of the fact that many of the more complicated CVs that we wish to compute have a functional form that can be thought of as

follows:

$$s = \sum_{i=1}^N f(\{X\}_i) \quad (13)$$

In other words, to calculate the CV one computes the same function,  $f$ , for  $N$  different sets,  $\{X\}_i$ , of atomic positions. Some more complicated examples of CVs of this form are the so-called secondary structure variables **ALPHARMSD**, **ANTIBETARMSD** and **PARABETARMSD** [22]. For **ALPHARMSD** one takes each set of  $M$  atoms in the protein that might together be able to form an alpha helix. For each of these sets of variables one then computes the RMSD distance between the instantaneous positions of the atoms and the positions of the atoms in an ideal alpha helix. Each of these RMSD distances is then transformed by a continuous switching function and these transformed values are then added together. The final value of the CV, that is calculated and printed using the input below, thus measures how many segments in the protein resemble an alpha helix.

```
MOLINFO STRUCTURE=helix.pdb
a: ALPHARMSD ...
  RESIDUES=all TYPE=DRMSD
  LESS_THAN={RATIONAL R_0=0.08 NN=8 MM=12}
...
PRINT ARG=a.less_than FILE=colvar STRIDE=10
```

The **ANTIBETARMSD** and **PARABETARMSD** commands do something very similar but the reference configuration in these cases are obviously an ideal anti-parallel beta sheet and an ideal parallel beta sheet respectively. For these three commands one could, in theory, calculate this quantity using multiple **DRMSD** commands and the **CUSTOM** command that was introduced earlier. It is obviously much easier though to use the input above, which determines the sets of atoms for which you have to determine RMSD distances from the template pdb structure directly. Having said that, however, it can be useful to try to write input files that use only the simple commands that were introduced in the previous sections in order to better understand what precisely is being calculated by these shortcuts.

This chapter would be overly long if we listed all the CVs that are available in PLUMED. If you are interested in these features we would thus recommend reading the literature. To get you started we have provided the following short, but far from exhaustive list, of references for a range of CVs, which can be calculated by PLUMED, and that includes the total energy of the system [23, 24], some of the components of the energy [25, 26], the dimer interaction energy [27], discrepancy measures [28], principal components [29, 30], path collective variables [31, 32], property maps [33], puckering variables [34, 35], Steinhart order parameters [17] and a number of experimental observables [36]. In addition, notice that PLUMED input files can be prepared using a graphical user interface [37] and that this graphical user interface also allows you to compute PLUMED CVs from within the VMD program [38].

### 3 Biasing collective variables

Section 2 has hopefully given you a sense of some of the quantities we might be interested in monitoring during the course of a simulation. If this were all that PLUMED could do, however, there would be no reason for it to be usable on the fly. After all, all of the quantities we have discussed can be calculated by post-processing the trajectory files that are output during the simulation. The reason that PLUMED runs on the fly is that it can also modify the Hamiltonian,  $H(\mathbf{q}, \mathbf{p})$ . These modifications introduce additional forces that must be incorporated when the underlying MD code integrates the equation of motion. Calculation of the PLUMED potential and the associated forces can, however, be separated from the calculation of the forces due to the underlying potential as the bias potential,  $V(\mathbf{q}, t)$ , that is calculated by PLUMED and that is a function of the atomic positions and possibly also time,  $t$ , is simply added to the Hamiltonian,  $H(\mathbf{q}, \mathbf{p})$ , that is calculated by the MD code. In other words, the modified Hamiltonian,  $H'(\mathbf{q}, \mathbf{p})$  is:

$$H'(\mathbf{q}, \mathbf{p}) = H(\mathbf{q}, \mathbf{p}) + V(\mathbf{q}, t) \quad (14)$$

#### 3.1 Reweighting

To understand reweighting we must return once again to the probability distribution that is sampled during a molecular dynamics trajectory. This distribution was first introduced in section 2.1 when we introduced the following proportionality:

$$P(\mathbf{q}, \mathbf{p}) \propto e^{-\frac{H(\mathbf{q}, \mathbf{p})}{k_B T}} \quad (15)$$

and explained how the quantity on the right-hand side of this equation is proportional to the probability of sampling any given vector of atomic positions,  $\mathbf{q}$ , and momenta,  $\mathbf{p}$ . Now suppose that we are not integrating the Hamiltonian  $H(\mathbf{q}, \mathbf{p})$  given above and that we are instead integrating a biased Hamiltonian  $H(\mathbf{q}, \mathbf{p}) + V(\mathbf{q})$ . The probability distribution,  $P'(\mathbf{q}, \mathbf{p})$ , that we will sample from when we integrate this biased Hamiltonian will be proportional to:

$$P'(\mathbf{q}, \mathbf{p}) \propto e^{-\frac{H(\mathbf{q}, \mathbf{p})}{k_B T}} e^{-\frac{V(\mathbf{q})}{k_B T}} \quad (16)$$

for the same reasons that the probability of sampling a particular configuration when the system is unbiased is given by equation 15. Notice, furthermore, that the right-hand side of equation 15 appears in the equation above and that we can thus rewrite this expression as:

$$P(\mathbf{q}, \mathbf{p}) \propto P'(\mathbf{q}, \mathbf{p}) e^{+\frac{V(\mathbf{q})}{k_B T}} \quad (17)$$

This equation relates the probabilities that we extract from biased simulations to the corresponding unbiased probabilities and, as we will see in the next section, it is thus the central equation for the analysis of biased simulations.

### 3.2 Extracting the free energy

We now have all the pieces we require and can thus finally discuss the topic that was first introduced in section 2; namely, how we analyze biased and unbiased MD simulations so as to extract the free energy as a function of a small number of collective variables. The first step in this process is to discretise the CV space and to introduce a set of probabilities,  $p_j$ , each of which tells us the probability that the CV value falls in a particular range or bin. These bins are set up so that all possible CV values fall within exactly one of the bins. Consequently, if we discretize the CV space,  $s(\mathbf{q}_i)$ , into  $M$  bins with centers at  $s_j$  the probability that  $t_1$  of the frames from our MD trajectory fall within the first of these bins, that  $t_2$  of frames fall within the second bin and so on is given by the following multinomial distribution:

$$P(t_1, t_2, \dots, t_M) \propto \prod_{j=1}^M p_j^{t_j} \quad (18)$$

where  $p_j$  is the probability that any given trajectory frame will fall into the  $j$ th bin. It is easy to extract the set of  $t_j$  values that appear in this expression from our trajectory. All we need to do is construct a histogram that tells us the number of times the trajectory visited each bin (see Note 4). Furthermore, once we have this information on the values of  $t_j$  we can perform a constrained optimization on equation 18 and thus determine the most likely values for the probabilities, the  $p_j$ s, given the particular set of  $t_j$  values that we observed during our trajectory. It is easier if we start this process of optimizing equation 18 by taking the logarithm of both sides of the equation. Doing so has no effect on the position of the minimum as the logarithm is a monotonically increasing function. Having taken the logarithm we then recall that the set of  $p_j$  values are probabilities and that as such they must satisfy  $\sum_{j=1}^M p_j = 1$ . We must therefore perform a constrained optimization using the method of Lagrange multipliers. The final function we need to optimize is thus:

$$\mathcal{L} = \sum_{j=1}^M t_j \ln p_j + \lambda \left( \sum_{j=1}^M p_j - 1 \right) \quad (19)$$

where  $\lambda$  is a Lagrange multiplier. When this expression is maximized we obtain (see Note 5)

$$p_k = \frac{t_k}{T} \quad (20)$$

where  $T = \sum_j t_j$  is the total number of frames in the trajectory. In other words, the most likely estimate for the probability that the CV will take a value within a particular range is just the fraction of time that the trajectory spent with the CV values in that particular range. By recalling the relationship between probability and free energy we can thus express the free energy for having the CV value in the  $k$ th range as:

$$F_k = -k_B T \ln t_k + C \quad (21)$$



Let us now suppose that we had taken the data from a biased trajectory rather than an unbiased trajectory and discuss how we would extract the histogram in this case. In other words, suppose that instead of integrating the Hamiltonian,  $H(\mathbf{p}, \mathbf{q})$ , directly we had integrated the biased Hamiltonian,  $H(\mathbf{p}, \mathbf{q}) + V(\mathbf{q})$ . As discussed in the previous section we know that the probability distribution that we sample configurations from when we do this biased simulation,  $P'(\mathbf{q}, \mathbf{p})$ , is related to the probability distribution that we would have sampled configurations from had we run our simulation without the bias,  $P(\mathbf{q}, \mathbf{p})$ , by equation 17. We therefore might expect that we can extract a free energy landscape for the unbiased simulation using the data from our biased simulations and a procedure much like that outlined above. The trick for doing this is to recognize that we can use equation 17 to write the multinomial distribution that we sample from when we run a *biased* simulations in terms of the elements of the probability distribution,  $p_j$ , that would have been sampled if we had run our simulation *without any bias* as shown below:

$$P'(t_1, t_2, \dots, t_M) \propto \prod_{j=1}^M (w_j p_j)^{t_j}, \quad (22)$$

where  $w_j = e^{-\frac{V(\mathbf{q}_j)}{k_B T}}$  and where  $V(\mathbf{q}_j)$  is the bias potential calculated for the  $j$ th frame. If we take the logarithm of this expression and impose the constraint that  $\sum_{j=1}^M (w_j p_j) = 1$  we thus find that the function to optimize in this case is:

$$\mathcal{L} = \sum_{j=1}^M t_j \ln w_j p_j + \lambda \left( \sum_{j=1}^M w_j p_j - 1 \right) \quad (23)$$

When this expression is maximized we obtain (see Note 6):

$$p_k \propto t_k e^{+\frac{V(\mathbf{q}_k)}{k_B T}} \quad (24)$$

By taking the negative logarithm of the above and by multiplying by  $k_B T$  we can thus extract the unbiased free energy surface from a biased simulation.

### 3.3 Biasing simulations with PLUMED

The need for numerous methods based on integrating modified Hamiltonians has been discussed at length in section II. We will thus limit the discussion of these methods here to providing a few sample PLUMED input files that instruct PLUMED to calculate various bias potentials. The following input instructs PLUMED to use a harmonic restraint on the distance between atoms 1 and 2 as a bias potential. The instantaneous values for the distance and the bias potential are then output to a file called colvar.

```
d1: DISTANCE ATOMS=1,2
rr: RESTRAINT ARG=d1 AT=0.2 KAPPA=10
PRINT ARG=d1,rr.bias FILE=colvar
```

This input encourages the system to sample configurations where the distance between atoms 1 and 2 is close to 2 Å. Many such inputs are typically used when performing umbrella sampling calculations with multiple restraints [39, 40].

Bias potentials do not have to be independent of time. There are, in fact, a whole class of steered MD methods in which a time dependent bias potential is used to force the system to change its conformation over the course of a simulation. A sample PLUMED input for such a calculation is given below.

```
MOLINFO MOLINFO STRUCTURE=helix.pdb
phi3: TORSION ATOMS=@phi-3
mr: MOVINGRESTRAINT ...
    ARG=phi3 STEP0=0 STEP1=10000 STEP2=20000
    KAPPA=100 AT0=-pi/3 AT1=pi/4 AT2=-pi/3
...
PRINT ARG=phi3,mr.* FILE=colvar
```

The above would be the input for a 20000 step steered MD simulation [41] in which the  $\phi$  angle in the third residue of protein is forced to change its value from  $-\frac{\pi}{3}$  radians to  $\frac{\pi}{4}$  radians before being forced to change the value of this torsional angle back to  $-\frac{\pi}{3}$  radians once more. This input instructs PLUMED to output the instantaneous value of the torsional angle, the bias and the work the bias has done on the system. Obviously, the user can make the path the system is forced to take through configuration space more complicated by using linear combinations of CVs and by changing the number of STEP and AT commands.

The final sample input below gives an example of how PLUMED can be used to perform the metadynamics simulations [42] that were discussed at length in Chapter 4.

```
phi: TORSION ATOMS=5,7,9,15
psi: TORSION ATOMS=7,9,15,17
metad: METAD ARG=phi,psi PACE=500 HEIGHT=1.2 SIGMA=0.35,0.35
PRINT STRIDE=10 ARG=phi,psi,metad.bias FILE=COLVAR
```

The input here is for a classic test case for these methods: calculating the free energy surface for alanine dipeptide as a function of the  $\phi$  and  $\psi$  backbone dihedral angles.

What you have hopefully noticed from the above input files is that PLUMED separates the calculation of the CVs from the calculation of the bias potential. Consequently, when using PLUMED we would normally express the bias as:

$$V(\mathbf{q}, t) = V(s_1(\mathbf{q}), s_2(\mathbf{q}), \dots, s_n(\mathbf{q}), t) \quad (25)$$

This introduces a set biased collective variables  $\mathbf{s}(\mathbf{q}) = \{s_1(\mathbf{q}), s_2(\mathbf{q}), \dots, s_n(\mathbf{q})\}$  that allow us to represent the dynamics that all the atoms undergo in some low-dimensional space. If we can calculate the partial derivatives of these collective variables with respect to the atomic positions we can then use the chain rule to calculate the forces on the atoms that result from the bias potential using:

$$f_m = -\frac{\partial V}{\partial q_m} = -\sum_{i=1}^n \frac{\partial V}{\partial s_i} \frac{\partial s_i}{\partial q_m} \quad (26)$$

Here the sum runs over the  $n$  collective variables that  $V(\mathbf{s}, t)$  is a function of.  $f_m$ , meanwhile, is the force acting on one of the three components of the position of a particular atom, and  $\frac{\partial s_i}{\partial x_m}$  is the derivative of the  $i$ th collective variable with respect to this particular component of the position of this same atom.

The key point to recognize is that the bias, which is some complicated function of the atomic positions, is expressed as a function of some set of simpler functions of the atomic positions. Consequently, as we saw in the discussion of CVs, we can thus build very complicated bias potentials by combining simpler pieces together in the PLUMED input file. For more information on the frequency with which CVs are calculated see Note 7.

As was the case for the CVs PLUMED contains implementations of other free energy methods that we do not have the space to discuss here. We will thus finish this section once more with a list of methods that are implemented in PLUMED together with relevant papers that describe them. This list of methods includes: many variants on metadynamics [43, 44, 45, 46, 47, 48, 49, 50, 11, 51, 52] including techniques that allow boundary conditions to be treated correctly in metadynamics simulations [53] and methods for extracting kinetic information from metadynamics simulations [54]. Then, in addition to metadynamics, we also have implementations of temperature-accelerated MD [55, 56], adaptive biasing force [57, 58, 59], and variationally enhanced sampling [60, 61]. Lastly, several bias potentials are included that allow one to force the ensemble averages that are extracted from the simulations to agree with data extracted from experiments. The implemented techniques include experiment directed simulations [62, 63], maximum entropy restraints [64], target metadynamics [65, 66, 67], and metainference [68].

### 3.4 Free Energies

To calculate the free energy as a function of the distance between two atoms using PLUMED and the technique described in section 3.2 one would use an input like the one shown below:

```
d1: DISTANCE ATOMS=1,2
h: HISTOGRAM ...
  ARG=d1 GRID_MIN=0 GRID_MAX=1.0
  GRID_BIN=200 KERNEL=DISCRETE STRIDE=10
...
f: CONVERT_TO_FES GRID=h TEMP=300
DUMPGRID GRID=f FILE=fes.dat
```

The input above creates a histogram with 201 bins that cover the range of distance values between 0 and 10 Å. This histogram is constructed using the distance calculated for every 10th trajectory frame. The free energy is then calculated from the accumulated histogram and output once the calculation has completed.

We can also use PLUMED and the methods described in section 3.2 to calculate the free energy landscape from a biased simulation. Below is a sample

input that performs this type of calculation:

```
phi: TORSION ATOMS=1,2,3,4
psi: TORSION ATOMS=5,6,7,8
EXTERNAL ARG=phi,psi FILE=bias_potential.dat

as: REWEIGHT_BIAS TEMP=300

hh: HISTOGRAM ...
    LOGWEIGHTS=as ARG=d STRIDE=10
    GRID_MIN=-pi,pi GRID_MAX=pi GRID_BIN=400
    KERNEL=DISCRETE
... HISTOGRAM

ff: CONVERT_TO_FES GRID=hh TEMP=300
DUMPGRID GRID=ff FILE=fes.dat
```

This is an input file for an enhanced sampling calculation that uses a bias that has been read from a file called `bias_potential.dat`. The bias potential that is read from this file is designed to act on the two torsional angles and to thus force them to sample configuration space more fully. The histogram is constructed using data from every 10th trajectory frame. The free energy is then calculated from the accumulated histogram and output once the calculation has completed. This is all very similar to what was done in the unbiased case. The major difference here, however, is that we have used the `REWEIGHT_BIAS` command and the `LOGWEIGHTS` keyword to ensure that the effect of the bias is discounted when the histogram is constructed. The free energy surface output is thus the free energy surface for the unbiased Hamiltonian.

It is important to notice that weighting each of the simulated snapshots with a factor of  $e^{+\frac{V(q)}{k_B T}}$  is correct when  $V$  only depends on the coordinates of the systems. If the bias potential,  $V$ , is time-dependent, more advanced procedures, that will not be discussed at length in this chapter, should be employed. Two particularly important enhanced sampling techniques that employ time dependent bias potential, and for which we provided sample inputs above, are:

- **Steered MD:** To extract free energies from such simulations one has to run multiple steered MD calculations and then analyze the ensemble of trajectories obtained using the Jarzynski equality [69].
- **Metadynamics:** Several procedures have been proposed that allow one to compute appropriate weighting factors [70, 48, 71].

## 4 Calculating error bars

The methods that we have introduced thus far for calculating ensemble averages and for calculating histograms are all based on the fact that we can calculate

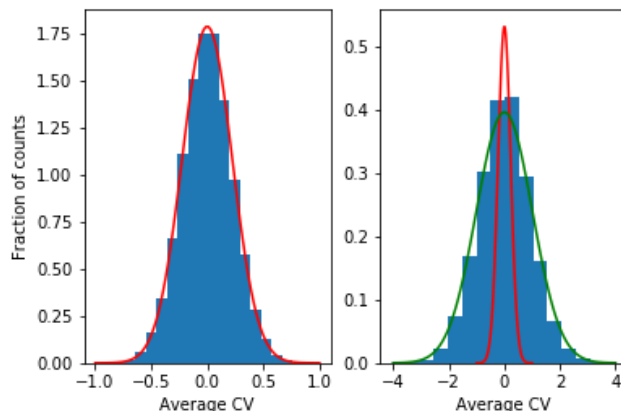


Figure 11: Histograms that are obtained by sampling independent and identically distributed random variables (left) and by sampling CV values from a trajectory (right). Each of the quantities that were used to calculate the histograms was constructed by taking an average of  $n = 20$  random variables. The central limit theorem would thus predict the distribution of these averages to be a Gaussian with a variance of  $\frac{\sigma^2}{n}$ , where  $\sigma^2$  is the sample variance. The red lines thus show what this analytical function looks like for the two data sets. For the independent random variables in the left panel it is clear that this is a good model for the sampled data. For the correlated data in the right panel, however, this model substantially underestimates the variance in the true distribution of averages. In fact, the green line shows a Gaussian that has the sample variance as its parameter and it is clear that this provides a much better model for the shape of this distribution.

an approximate value for the true ensemble average by taking a large number of samples from the underlying distribution. In other words, in all these methods we are calculating a sample mean and assuming that it provides a reasonable estimate for the true population mean. What has thus far been missing from this discussion, however, is how we get an estimate for the error bar on the sample mean. Any sample mean, after all, is itself a random variable and thus has an underlying distribution, which we should attempt to characterize. The difficulty in doing so when our data is extracted from an MD trajectory is, however, that there are strong and measurable (see Note 8) correlations between the configurations sampled in the frames of the trajectory that are adjacent. As we will see in the remainder of this section, however, we can take these correlations into account when we calculate the variance and can thus calculate error bars.

Figure 11 shows why these correlations represent such a problem. To construct the histogram shown in the panel on the left we generated 50000 sets of  $n = 20$  normal random variables. We then calculated a sample average for

each of these sets,  $\mu_i$ , and a sample mean,  $\langle X \rangle$ , and sample variance,  $\sigma^2$ , for all 1000000 points in the data set. This sample variance was calculated using:

$$\sigma^2 = \frac{1}{N-1} \sum_{i=1}^N (X_i - \langle X \rangle)^2 \quad (27)$$

The reason for doing this is that the central limit theorem tells us that the distribution of  $\mu_i$  values should be a Gaussian centered on the sample mean with a variance of  $\frac{\sigma^2}{n}$ . As you can see from the left panel this is a good model for the distribution of the sample means that was obtained by averaging the independent random variables. The right panel shows, however, that this is not a good model when the data is taken from an MD or Monte Carlo trajectory. As discussed in Note 9 to construct this figure sample means for 50000 sets of  $n = 20$  correlated random variables were calculated once more as well as the sample mean and variance for the full set of 1000000 variables. The blue histogram shows the distribution of the sample means. As you can see the correlation between the values of the CVs ensures that this distribution is substantially wider than the central limit theorem (red line) would predict it to be. In other words, for correlated data  $\frac{\sigma^2}{n}$  is not a sensible estimate for the variance of the sample mean.

This problem can be resolved using the so-called block averaging technique [72]. In essence this technique involves splitting the trajectory into a set of  $N$  blocks that are all of equal length.  $N$  separate sample means,  $\mu_i$ , are then calculated from each of these blocks of data. Obviously, the mean for these  $N$  quantities that we calculate using:

$$\mu = \frac{1}{N} \sum_{i=1}^N \mu_i \quad (28)$$

is equal to the mean we would have obtained had we averaged every frame in our trajectory. The advantage of block averaging in this way is, however, that, if each  $\mu_i$  is calculated over a long enough block of trajectory, the values of the  $\mu_i$ s are no longer correlated. We can thus calculate an error bar on our estimate of  $\mu$  that we would calculate using the formula above as:

$$\epsilon = \Phi \left( \frac{p_c + 1}{2} \right) \sqrt{\frac{1}{N(N-1)} \sum_{i=1}^N (\mu_i - \mu)^2} \quad (29)$$

where  $\Phi$  is the cumulative probability distribution function for the standard normal random variable and where  $p_c$  is the level of statistical confidence we want our error bar to represent.

Figure 12 shows the outcome of applying this technique on data obtained by sampling independent and identically distributed random variables and on data obtained by calculating the value of a CV over the course of a trajectory. The mean and error bars are shown as a function of the sizes of the blocks over

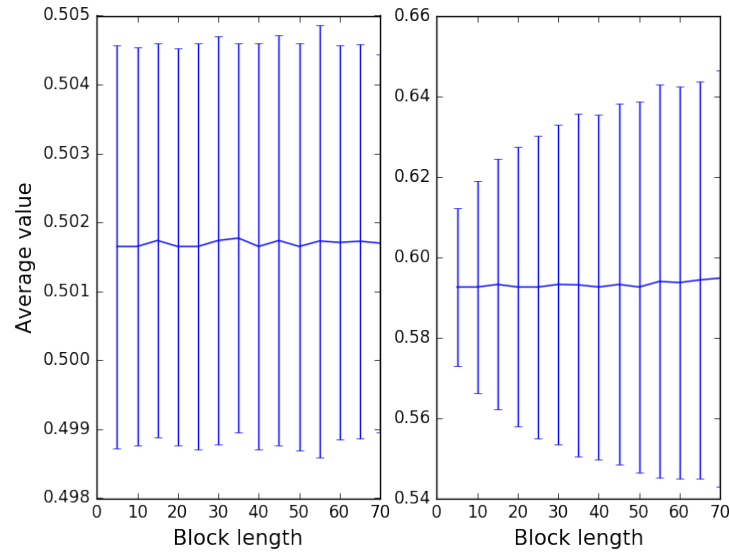


Figure 12: Figure showing how the error bars calculated using the block averaging technique depend on the lengths of the blocks used for independent and identically distributed random variables (left) and for trajectory data (right). As you can see, the size of the error bar is largely independent of the block size when the data points being averaged are samples of a random variable. The correlations between the CV values we obtain in a trajectory, however, ensure that the error bar is underestimated when the block size used is small. Consequently, as the length of time over which the block averages are taken is increased the error bar increases in size until eventually it reaches a plateau.

which the  $\mu_i$  averages were calculated. In both cases you can see that the final value of  $\mu$  does not depend on the size of the blocks. The small differences in this value are due to the finite precision algebra the computer uses. For the correlated data from the MD trajectory, however, the size of the error bar does depend on the length of the blocks. In particular, the error bar is smaller when the block sizes are small as the  $\mu_i$  values are still correlated. As the block sizes get larger, however, the  $\mu_i$  values become decorrelated and the size of the error bar thus plateaus to a near constant value.

A similar technique can be used to calculate error bars for histograms and free energy surfaces. For the weighted histograms discussed at the end of the section 3.2 the probability,  $g_j^{(i)}$ , for the  $j$ th bin of the histogram is computed using equation 24, and the data from the  $i$ th block of trajectory data. This procedure gives multiple estimates for the probability of each bin, which can be combined to give a single set of sample means and confidence limits using the following or other similar expressions:

$$\langle g_j \rangle = \frac{1}{\sum_{i=1}^N w_i} \sum_{i=1}^N g_j^{(i)} \quad (30)$$

$$\epsilon_j = \Phi\left(\frac{p_c + 1}{2}\right) \sqrt{\frac{W}{N(W - \frac{S}{W})} \sum_{i=1}^N w_i * (g_j^{(i)} - \langle g_j \rangle)^2} \quad (31)$$

Here  $w_i$  is used to indicate the sum of the weights of all the configurations in each of the blocks of data.  $W$  and  $S$  are then the sum of these  $w_i$  values and the sum of the squares of these values respectively.

Example python code for taking the set of histogram files output by PLUMED and for calculating the means and the error bars using the formulas above can be found on the PLUMED website. Figure 13 shows how the average size of the errors on the estimate of the free energy obtained from a metadynamics simulation changes as the length of the blocks is increased (see Note 10). Similarly to what we observed in the right panel of figure 12 when small blocks are used the error bar on the free energy is underestimated and the error bars are thus too small. When sufficiently large blocks are used, however, the size of the error reaches a plateau value.

## 4.1 Block analysis with PLUMED

We can use PLUMED to perform the block averaging described in section 4 and can thus obtain error bars on the histograms we extract from biased or unbiased trajectory data. The following PLUMED input can be used to calculate an independent histogram from each block of 1,000,000 MD steps in the trajectory. The simulation here is the same biased calculation that was described in section 3 and the histograms are output to a set of files called analysis.0.histo.dat, analysis.1.histo.dat, ... and histo.dat.



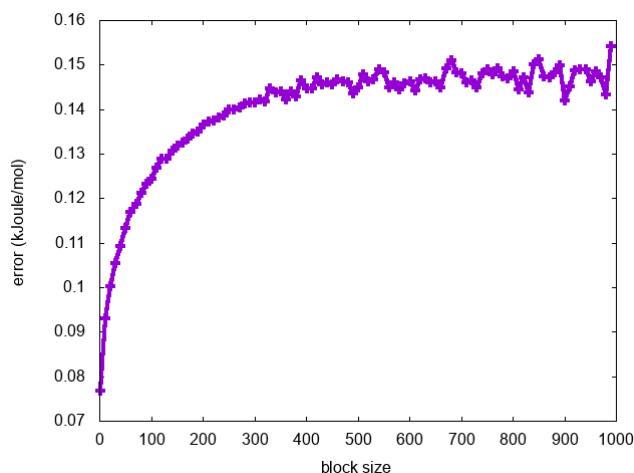


Figure 13: Figure showing how the average error on the estimate of a free energy obtained for a metadynamics simulation of alanine dipeptide depends on the sizes of the blocks over which the individual histograms were calculated. As with the example involving the ensemble average in figure 12 the error is underestimated when small blocks are used. When larger blocks of data are used, however, the value of the error plateaus at a constant value.

```
phi: TORSION ATOMS=1,2,3,4
psi: TORSION ATOMS=5,6,7,8
EXTERNAL ARG=phi,psi FILE=bias_potential.dat
```

```
as: REWEIGHT_BIAS TEMP=300
```

```
hh: HISTOGRAM ...
    LOGWEIGHTS=as ARG=d STRIDE=10
    GRID_MIN=-pi,pi GRID_MAX=pi GRID_BIN=400
    KERNEL=DISCRETE CLEAR=1000000
... HISTOGRAM
```

```
DUMPGRID GRID=hh FILE=histo.dat STRIDE=1000000
```

The probabilities from all these histograms can then be combined using the equations in section 4 although as mentioned in the previous section you can find scripts on the PLUMED website that will compute these quantities for you.

## 5 Multiple replica techniques

Thus far we have considered cases where a single simulation is performed and then analyzed. There are, however, situations where one wants to simultane-

ously run multiple simulation replicas. Sometimes these replicas are run under equivalent conditions and the only differences between them are the initial configurations. If this is the case one can simply concatenate all the trajectories and analyze them using the techniques that have been discussed in the previous sections. Alternatively, if the simulations are short and thus not expected to reach equilibrium, more statistically reliable results may well be obtained if they are analyzed using the Markov state models that have been discussed in chapter IV. Oftentimes, however, the various replicas are run using different conditions (e.g., replicas have different temperatures or different biasing potentials). Furthermore, in this case the replicas will often communicate with each other using a replica-exchange procedure. In this section we will thus discuss how such multiple-replica simulations operate and how they should be analyzed. The analysis described in what follows can be applied when replicas are simulated independently of each other or when they are simulated in a manner that permits exchanges between replicas. For best results we would always suggest allowing exchanges between replicas as this increases the ergodicity of the simulation.

Replica-exchange molecular dynamics provides a rigorous theoretical framework that can be used to swap coordinates between trajectories that are calculated simultaneously but performed using different conditions (see Chapter II). Perhaps the most commonly used of these techniques is parallel tempering [73], which uses a different temperature in each of the simulated replicas. The replicas can also experience different simulation biases, however, so, because PLUMED is designed to bias MD simulations and because there are a number of techniques in the literature that are based on this principle (see, e.g., [74, 46, 75, 11]), this will be the focus in this section. Within many of these replica-exchange techniques the Metropolis Monte Carlo scheme is employed with proposed moves that involve swapping the coordinates of replica  $i$  ( $\mathbf{q}^{(i)}$ ) with the coordinates of replica  $j$  ( $\mathbf{q}^{(j)}$ ) and an acceptance probability that is given by:

$$\alpha = \min \left( 1, \frac{P^{(i)}(\mathbf{q}^{(j)})P^{(j)}(\mathbf{q}^{(i)})}{P^{(i)}(\mathbf{q}^{(i)})P^{(j)}(\mathbf{q}^{(j)})} \right) \quad (32)$$

Here  $P^{(i)}(\mathbf{q})$  is the probability that the set of coordinates  $\mathbf{q}$  will be observed under the conditions experienced by the  $i$ th replica. Assuming that the only difference between the two replicas is the bias potential experienced and that the bias potentials for replica  $i$  is  $V^{(i)}(\mathbf{q})$  allows us to use equation 16 to rewrite this probability using an expression that only depends on the bias potentials:

$$\alpha = \min \left( 1, \exp \left( \frac{-V^{(i)}(\mathbf{q}^{(j)}) - V^{(j)}(\mathbf{q}^{(i)}) + V^{(i)}(\mathbf{q}^{(i)}) + V^{(j)}(\mathbf{q}^{(j)})}{k_B T} \right) \right) \quad (33)$$

When coordinate swaps are accepted or rejected using these criteria we ensure that the system remains at equilibrium. Consequently, any averages that we obtain are in principle indistinguishable from those that would have been obtained if the simulations had been performed independently. The advantage, however, is that the coordinate swaps usually systematically increase the number of slow

transitions that are sampled. We can thus use these methods to obtain more statistically robust averages from shorter (parallel) simulations.

## 5.1 The weighted histogram method

The multiple replica simulations that were introduced above provide us with a powerful set of tools that allow us rapidly explore a wide region of configuration space and to take advantage of parallel computing facilities. In this final section we will discuss how the trajectories we obtain from these simulations are analyzed. Much like the analysis that was discussed in section 3.2 the aim here is to extract the free energy as a function of a CV or set of CVs. Furthermore, when analyzing these multiple replica simulations we are again going to calculate a histogram and then exploit the maximum likelihood technique to convert this to a free energy. At odds with section 3.2, however, we now have multiple trajectories, each of which was obtained by integrating a different biased Hamiltonian. We thus calculate the probability of observing this particular set of configurations during the  $N$  trajectories that we ran using the product of multinomial distributions shown below:

$$P(\mathbf{T}) \propto \prod_{j=1}^M \prod_{k=1}^N (c_k w_{kj} p_j)^{t_{kj}} \quad (34)$$

In this expression the second product runs over the biases that were used when calculating the  $N$  trajectories. The first product then runs over the  $M$  bins in our histogram. The  $p_j$  variable that is inside this product is the quantity we wish to extract; namely, the unbiased probability of having a set of CV values that lie within the range for the  $j$ th bin. The quantity that we can easily extract from our simulations,  $t_{kj}$ , however, measures the number of frames from trajectory  $k$  that are inside the  $j$ th bin. To interpret this quantity we must consider the bias that acts on each of the replicas so the  $w_{kj}$  term is introduced. This quantity is calculated using equation 17 and is essentially the factor that we have to multiply the unbiased probability of being in the bin by in order to get the probability that we would be inside this same bin in the  $k$ th of our biased simulations. Obviously, these  $w_{kj}$  values depend on the value that the CVs take and also on the particular trajectory that we are investigating all of which, remember, have different simulation biases. Finally,  $c_k$  is a free parameter that ensures that, for each  $k$ , the biased probability is normalized.

We can use equation 34 to find a set of values for  $p_j$  that maximizes the likelihood of observing the trajectory. This constrained optimization must be performed using a set of Lagrange multipliers,  $\lambda_k$ , that ensure that each of the biased probability distributions are normalized, that is  $\sum_j c_k w_{kj} p_j = 1$ . Furthermore, much as in section 3.2, the problem is made easier if equation 34 is replaced by its logarithm.

$$\mathcal{L} = \sum_{j=1}^M \sum_{k=1}^N t_{kj} \ln c_k w_{kj} p_j + \sum_k \lambda_k \left( \sum_{j=1}^N c_k w_{kj} p_j - 1 \right) \quad (35)$$

After some manipulations (see Note 11), the following equations emerge:

$$p_j \propto \frac{\sum_{k=1}^N t_{kj}}{\sum_k c_k w_{kj}} \quad (36)$$

$$c_k = \frac{1}{\sum_{j=1}^M w_{kj} p_j} \quad (37)$$

Equations 36 and 37 can be solved by computing the  $p_j$  values with an initial guess of the  $c_k$  values and by then refining these values using the  $c_k$  values that are obtained by inserting the  $p_j$  values obtained into equation 37. Usually the  $c_k$  and  $p_j$  values become self-consistent after a few rounds of this iterative algorithm.

When writing scripts to do this form of analysis it is worth noting that only  $\sum_k t_{kj}$ , which is the total number of configurations from *all* the replicas that enter the  $j$ th bin, enters equations 36 and 37. There is thus no need to record which replica generated each of the frames and one can thus simply gather the trajectories from all the replicas together at the outset.

The fact that we can simply gather the trajectories from all our replicas before performing the weighted histogram analysis that has been described in the previous paragraph is also evident when we consider equation 36. This expression tells us that we can calculate  $p_j$  by constructing a weighted histogram from the concatenated trajectory and that the weight for the  $n$ th frame will be:

$$\tilde{w}_n = \frac{1}{\sum_k c_k w_{kn}} \quad (38)$$

where the sum runs over the replicas and where  $w_{kn}$  is the factor the  $k$ th replica has to be reweighted by in order to recover the unbiased probability for configuration  $n$ . Notice, that we can also use the concatenated trajectory when extracting values for  $c_k$  using equation 37 as the sum over bins in the denominator can be replaced with a sum over the concatenated trajectory. These observations are important as they are the basis of the binless formulation of the weighted histogram technique (WHAM) that is implemented within PLUMED and that has been variously proposed by a number of different authors [76, 77, 78]

## 5.2 WHAM analysis with PLUMED

Section 5 discussed simulations in which multiple replicas are run in parallel and in which exchanges are attempted between replicas. To run these types of calculations you need an MD code that can manage the communication between these replicas. GROMACS is able to run simulations in this way using a command line syntax similar to the one shown below, which tells GROMACS to run 3 replicas in parallel and to attempt coordinate exchanges every 1000 MD steps:

```
> mpirun -np 3 gmx_mpi mdrun -multi 3 -plumed plumed.dat -replex 1000
```

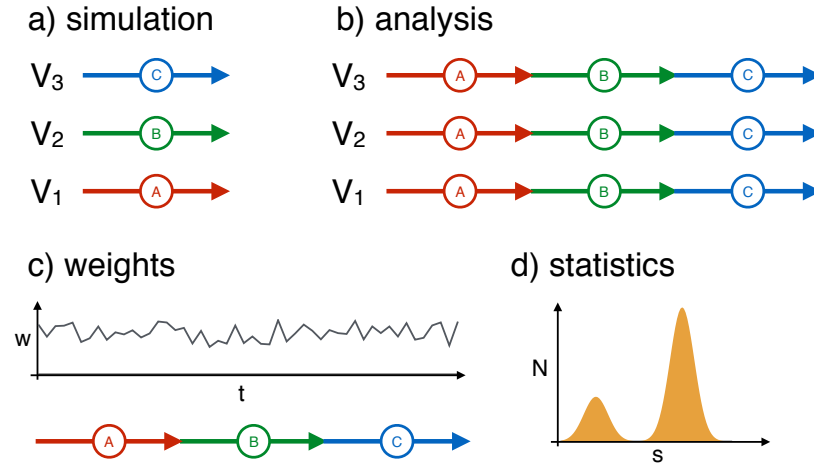


Figure 14: Schematic representation of a WHAM analysis. a) Multiple simulations are performed using different bias potentials. These simulations can be done separately, but we would always recommend using a replica-exchange procedure. b) Trajectories are concatenated for subsequent analysis. The first step of this analysis is to compute the value of all the biasing potentials on all the snapshots in the concatenated trajectory. Notice that in principle it is not necessary to concatenate the trajectory as the order of the accumulated snapshots is irrelevant. It is critical, however, to compute all the bias potentials for all the frames in all the trajectories as only then can you solve equations 36 and 37. c) Once these equations are solved a weight for each snapshot in the concatenated trajectory is obtained. d) These weights are then used to, for example, reconstruct the unbiased probability along some *a posteriori* chosen CV.

For MD codes that can handle multiple replicas PLUMED provides a convenient syntax for having different biases on the various different replicas. As an example the PLUMED input below assumes that three replicas are being run in parallel and that these three replicas differ only in the position of the center of the harmonic restraint:

```
d: DISTANCE ATOMS=10,20
RESTRAINT ARG=d AT=@replicas:1.0,1.1,1.2 KAPPA=1.0
```

Obviously, the centers of the harmonic restraint in the three simulations are at 1.0, 1.1 and 1.2 respectively. The CV on which this restraint acts and the strength of the restraint are, however, the same in all three replicas.

Once the multiple replica simulation has run, it must be analyzed. As discussed in section 5.1 the WHAM technique provides a good method for doing this analysis. To do WHAM using PLUMED you must first *concatenate* the trajectories from the various replicas. The exact way this will be done will depend on the format of the trajectory file. If the format is a plain text `.gro` file, file concatenation may be sufficient. For other file types, however, it may be necessary to use the specific tools that are provided by the MD engine. Regardless of these details, however, once a single concatenated trajectory is available it can be analyzed using PLUMED with an input file like the one shown below:

```
d: DISTANCE ATOMS=10,20
RESTRAINT ARG=d AT=@replicas:1.0,1.1,1.2 KAPPA=1.0
```

```
as: REWEIGHT_WHAM TEMP=300 ARG=d.bias
```

```
cc: COLLECT_FRAMES ARG=d STRIDE=1 LOGWEIGHTS=aa
```

```
hh: HISTOGRAM ...
    ARG=cc.d STRIDE=10
    GRID_MIN=0 GRID_MAX=10 GRID_BIN=100
    KERNEL=DISCRETE
... HISTOGRAM
```

```
ff: CONVERT_TO_FES GRID=hh TEMP=300
DUMPGRID GRID=ff FILE=fes.dat
```

The first part of this input will be basically identical to the input used for the biased calculations. The rest, meanwhile, is very similar to the input for reweighting that was discussed in section 3.4 although there are some important differences, the most important of which is the action `COLLECT_FRAMES` between the action `REWEIGHT_WHAM` and the action `HISTOGRAM`. The reason for these differences is that the weights associated to each snapshot by the action `REWEIGHT_WHAM` are only available once whole trajectory has been processed. This special syntax thus instructs the action `HISTOGRAM` to wait until the end of the trajectory and to only then retrieve the weights and construct the histogram. Notice also that equations 36 and 37 are solved in the `REWEIGHT_WHAM`

command and that this command can thus accept additional arguments in order to fine tune the tolerance with which these equations are solved. As a final note it is important to remember that there are multiple replicas when running the WHAM calculation using the input above, which means that PLUMED **driver** has to be called using a command like that shown below:

```
> mpirun -np 3 plumed driver --multi 3 --plumed plumed.dat
```

We will now see how to use this syntax can be used to compute the energy landscape of an adenosine in water. The simulation reported here was done using the setup and parameters described in Ref. [64]. Consequently, a simulation with 16 replicas was run using the command below:

```
> mpirun -np 16 gmx_mpi mdrun -multi 16 -plumed plumed.dat -replex 1000
```

with the following PLUMED input file

```
MOLINFO STRUCTURE=adenosine.pdb
chi: TORSION ATOMS=@chi-1
#
# Impose an umbrella potential on chi
# with a spring constant of 80 kJoule/mol
# and centered in chi=AT
#
r: RESTRAINT ...
  ARG=chi KAPPA=80.0
  AT=@replicas:{
    0*pi/8 1*pi/8 2*pi/8 3*pi/8
    4*pi/8 5*pi/8 6*pi/8 7*pi/8
    8*pi/8 9*pi/8 10*pi/8 11*pi/8
    12*pi/8 13*pi/8 14*pi/8 15*pi/8
  }
...
```

The numbers listed after the **@replicas** instruction are basically 16 equally-spaced values between 0 and  $2\pi$ . Furthermore, in setting up this grid of restraints we have exploited the fact that PLUMED can perform simple algebraic calculations when interpreting its input.

The restraints on the 16 replicas that will be simulated by executing the command above ensure that all the possible values for the  $\chi$  glycosidic torsion of this nucleoside, including the unfavorable ones that correspond to free-energy barriers, will be explored during these simulations. Furthermore, once the simulation has completed we can concatenate all the trajectories produced by GROMACS (called **traj0.xtc**, **traj1.xtc**, ... **traj15.xtc**) into a single long trajectory called **traj-all.xtc**. This trajectory can then be analyzed using the following command

```
> mpirun -np 16 plumed driver --multi 16 --plumed plumed.dat \
  --ixtc traj-all.xtc
```

and a PLUMED input file that is very similar to the one described above. Instead of writing a new file from scratch, it is often convenient to include the file that was used when the simulation was run in the analysis file by using the command `INCLUDE`. By doing so we can thus write an analysis file, called `plumed.wham.dat`, that reads as follows:

```
INCLUDE FILE=plumed.dat
# also compute the puckering of the sugar:
puck: PUCKERING ATOMS=@sugar-1
bias: REWEIGHT_WHAM TEMP=300 ARG=r.bias

cc1: COLLECT_FRAMES ARG=chi,puck.Zx LOGWEIGHTS=bias
h1: HISTOGRAM ...
    ARG=cc1.chi
    GRID_MIN=-pi
    GRID_MAX=pi
    GRID_BIN=100
    BANDWIDTH=0.1
...
fes1: CONVERT_TO_FES TEMP=300 GRID=h1
DUMPGRID GRID=fes1 FILE=fes1.dat

cc2: COLLECT_FRAMES ARG=chi,puck.Zx LOGWEIGHTS=bias
h2: HISTOGRAM ...
    ARG=cc2.chi,cc2.puck.Zx
    GRID_MIN=-pi,-pi
    GRID_MAX=pi,pi
    GRID_BIN=100,100
    BANDWIDTH=0.1,0.1
...
fes2: CONVERT_TO_FES TEMP=300 GRID=h2
DUMPGRID GRID=fes2 FILE=fes2.dat
```

For RNA, it is common to analyze the conformation of the sugar using the `Zx` component that is defined in [35]. To calculate this CV using PLUMED you use the `puck.Zx` component the `PUCKERING` command. Notice that this CV is used in the above input because, in addition to computing the free energy as a function of the  $\chi$  torsion it also computes a second free energy surface that depends on both  $\chi$  and the puckering conformation of the sugar. This second variable was not biased, but can also be analyzed at this latter stage. The two free energy surfaces that are extracted when the above analysis is performed on the trajectory are shown in figure 15. In addition, this figure also shows the result of using a `HISTOGRAM` command *without* taking the weighting factors into account.

Notice that the WHAM procedure discussed here can be used to remove any arbitrary bias that has been added to a simulation. For instance, it can be used to analyzed bias-exchange metadynamics simulations [46], which is a method



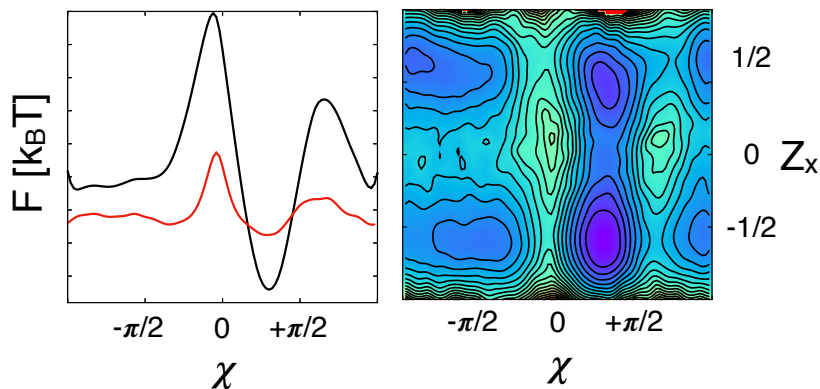


Figure 15: Results from a replica-exchange umbrella sampling simulation performed on adenosine in water. As explained in the text restraints are applied on the glycosidic torsion  $\chi$  in these simulations. (a) The free energy as a function of  $\chi$  computed using the WHAM equation is shown in black. This profile is in agreement with that reported in reference [64]. The free energy that one would have obtained if one had (erroneously) collected the histogram from the replica exchange simulation without performing any reweighting is also shown (red line). This profile is significantly flatter because the system is forced by the restraints to explore the entire CV span. (b) Two-dimensional profile showing the free energy as a function of the biased  $\chi$  torsion and the  $Z_x$  puckering variable, that reports on the sugar conformation. Even though this latter variable was not biased, the two-dimensional profile can be correctly reconstructed as transitions between the two typical conformations are accessible on the simulation timescale.

that involves a replica-exchange scheme in which each replica experiences a different metadynamics bias. In these types of calculations different replicas can use different PLUMED input files and can even use a different number of biasing potentials. For example, in references [21, 79], a similar procedure was even used to reweight bias-exchange metadynamics simulations where, in addition to the metadynamics potential, each replica was subject to a different restraint. In all these cases, however, the replicas always share the same simulation parameters. In other words, the only differences are in the bias that is applied by PLUMED.

## 6 Perspectives

This chapter has introduced the PLUMED code, explained some of the theory behind the methods that are implemented in this code and given some practical examples that show how PLUMED can be used. Space constraints mean that we cannot describe everything that PLUMED can do, so we have instead chosen to focus on some of the issues that come up most frequently on the code’s mailing list. The discussion in the previous sections has thus included information on the difficulties that can arise when treating periodic boundary conditions and a discussion on the proper treatment of statistical errors and the analysis of multiple replica simulations. In what follows we will finish by giving a brief perspective on how we believe that PLUMED will evolve in the near future.

Our principal reason for developing PLUMED was to provide functionality for performing enhanced sampling calculations. There were two interrelated reasons for writing this as a separate piece of software. The first of these was that we wanted an inter-operable implementation that could work with multiple MD codes. Different MD codes contain different functionalities and as such the particular engine that you will work with on any given project will depend on the particular system under study. At the same time, however, the efficacy of any enhanced sampling method depends strongly on the degree to which the biased CVs separate the metastable basins and transition states in the energy landscape. Consequently, PLUMED needed to be a large, stand-alone code not simply because some of the biasing methods that are implemented within are rather sophisticated but also because it contains implementations for many of the, in some cases very-complicated, CVs that have been used in these types of calculations. In fact, to our knowledge there is no other code that contains implementations of as many different CVs. Having said that, however, a number of recently developed CVs are not yet included and will thus be implemented in the near future. We are particularly interested in some of the ideas from machine learning [80, 81, 82, 83, 84, 85] that are entering this field and are actively investigating how such methods could be implemented within the PLUMED framework. Furthermore, we have in the last couple of years been working with the developers of Open Path Sampling [86] in order to help them access the structural features that can be calculated using PLUMED so that methods such as transition path sampling and forward flux can be performed on-the-fly using a broad range of CVs without any need to perform a posteriori analysis.

Another issue that we need to work on in the future is the performance of the code. As PLUMED is designed to compute and bias CVs its design is rather different from some of the other pieces of software that are used to drive MD simulations. One big potential performance bottleneck comes about because PLUMED needs access to a subset of the simulated atoms during *every time step*. This requirement can slow down the simulations particularly when the CVs are computed from the positions of a large number of atoms. In the future we will thus work in order to decrease the computational cost of the operation that transfers the atomic positions from the MD code to PLUMED. One easy way to reduce this cost, which will be available in version 2.5, is to allow expensive CVs that are computed within the underlying MD code to be transferred to PLUMED. Computing the CV in within the MD code would allow you take advantage of the data structures in the MD code. Furthermore, as an added benefit, one could also take advantage of this feature when using CVs that are based on features that are difficult to transfer to PLUMED such as partial charges or other functions of the electronic structure. It is still to be seen if any these features will find a practical application, however.

One thing that will not change a great deal in the future is PLUMED’s API. The simple and well-documented interface between PLUMED and the various MD codes that call it is one of the codes great strengths so any future change will be made in a way that ensures backward compatibility. The reason this interface is so important is that it is this that allows such a large number of MD codes to call PLUMED. In addition, some developers have incorporated the PLUMED API within their code base so that users can download their codes and immediately use PLUMED without performing any sort of patching procedure. This model of having the calls to PLUMED within the MD codes is something that we would like to use more widely in future and to facilitate this and to keep pace with the growing popularity of python within the molecular simulation community we have recently provided an interface to the PLUMED API that makes the PLUMED routines callable from python.

The most important changes to PLUMED that we foresee, however, will be the features contributed to the code by independent groups. It is clear from the github pages of PLUMED that a number of forks of the code are now being actively developed. Furthermore, some of the features that have been developed in these forks have been already contributed back into the main version of the code. We are convinced that transforming this project in a community effort will be the best way to keep it lively and up to date. In fact, inviting contributions from the whole simulation community is the only way to ensure that the code contains the most exciting recent methodological developments from the field.

## 7 Notes

1. The first version of PLUMED was released in 2009 [1]. A complete rewrite, that made the code easier to maintain and extend, was published in 2014 [2]. This second version changed the inner structure of the code and also

changed the syntax for the input file so it is this second version of the code that will be the focus of this chapter.

2. When a reference pdb file is provided using the MOLINFO command numerous shortcuts can be employed when calculating backbone torsional angles in proteins and nucleic acids. For example the input below instructs PLUMED to calculate and print the  $\phi$  and  $\psi$  angles in the 3rd and 9th residue of a protein.

```
MOLINFO STRUCTURE=helix.pdb
phi3: TORSION ATOMS=@phi-3
psi3: TORSION ATOMS=@psi-3
phi9: TORSION ATOMS=@phi-9
psi9: TORSION ATOMS=@psi-9
PRINT ARG=phi3,psi3,phi9,psi9 FILE=colvar STRIDE=10
```

A number of other convenient shortcuts are explained in the PLUMED manual.

3. If for some reason one only wishes to only disregard translation of the center of mass, that is to say if one wishes to include any displacements that come about because of rotation of the reference frame, when computing the RMSD one can replace TYPE=OPTIMAL with TYPE=SIMPLE. In addition, one can use one set of atoms to calculate the rototranslation operation that minimizes the RMSD and then use a different set of atoms to compute the final RMSD. In fact this technique is commonly used to track the position of a ligand in the reference frame of a protein.
4. The amount of time that the system spent in the bin  $j$  can be computed as follows:

$$t_j = \sum_{i=1}^T H\left(\frac{|s(\mathbf{q}_i) - s_j|}{w}\right) \quad \text{where} \quad H(x) = \begin{cases} 1 & \text{if } x < 1 \\ 0 & \text{otherwise} \end{cases} \quad (39)$$

where  $w$  is the width of each bin.

5. When we differentiate  $\mathcal{L}$  in Eq. 19 with respect to  $p_k$  we find that at the constrained optimum:

$$\frac{\partial \mathcal{L}}{\partial p_k} = \frac{t_k}{p_k} + \lambda p_k = 0 \quad \rightarrow \quad p_k = -\frac{t_k}{\lambda} \quad (40)$$

We know, however, that:

$$\sum_{j=1}^M p_j = 1 \quad \rightarrow \quad \lambda = -\sum_{j=1}^M t_j \quad (41)$$

If we add together the number of trajectory frames in each of bins, however, we get the total number of trajectory frames,  $T$ .  $\lambda$  is thus equal to  $-T$  and thus the most likely value of  $p_k$  is simply:

$$p_k = \frac{t_k}{T} \quad (42)$$

6. When we differentiate  $\mathcal{L}$  in Eq. 23 with respect to  $p_k$  and set its derivatives to zero we obtain:

$$\frac{\partial \mathcal{L}}{\partial p_k} = \frac{t_j}{p_k} + w_k \lambda = 0 \quad \rightarrow \quad p_k = -\frac{(w_k)^{-1} t_k}{\lambda} \quad (43)$$

We can then recall our constraint; namely  $\sum_{j=1}^M w_j p_j = 1$ . Notice that by enforcing this constraint we are not doing anything to ensure that the unbiased distribution,  $p_j$ , is normalized. The constraint instead ensures that the biased distribution  $w_j p_j$  is normalized. This is important as it is this probability distribution that is used when we compute the total probability of observing the trajectory. It is, in fact, not at all necessary for the unbiased probability distribution,  $p_j$ , to be normalized and to be clear the unbiased distribution that emerges when this form of analysis is performed will be unnormalized as we will obtain  $\lambda = -T$  and hence  $p_k = \frac{(w_k)^{-1} t_k}{T}$ . The final result is thus:

$$p_k \propto (w_k)^{-1} t_k \quad (44)$$

7. The **STRIDE** keyword takes a default value of one and tells you how frequently a PLUMED command is executed. When it is used in combination with the **PRINT** command, it thus controls the frequency with which the CVs are printed. In addition, PLUMED automatically knows that these CVs should only ever be calculated when they are printed. The **STRIDE** keyword can also be used with commands that bias CVs such as **RESTRAINT** and **METAD**, however. In this context the command tells PLUMED that the bias, and the biased CVs, should be computed with a frequency as part of a multiple-time step scheme [87, 88]. By using these schemes you can speed up calculations especially when expensive CVs are used. You should, however, only ever use moderate values for **STRIDE** in this case – typically, something between 1 and 5.
8. We can measure the degree of correlation within a time series,  $X_t$ , of random variables with expectation  $\langle X \rangle$  and variance  $\langle (\delta X)^2 \rangle$  by measuring the autocorrelation function:

$$R(\tau) = \frac{\langle (X_t - \langle X \rangle)(X_{t+\tau} - \langle X \rangle) \rangle}{\langle (\delta X)^2 \rangle} \quad (45)$$

The value of this function at  $\tau$  gives a measure of the average degree of correlation between each pairs of random variables that were measured  $\tau$

time units apart. If the random variables are all independent and identically distributed this function will decay to zero for all  $t > 0$ . If the autocorrelation function is calculated for a time series of CV values taken from a trajectory the function will not decay immediately, however, as the system will most likely not diffuse from one edge of CV space to the other during a single timestep. The CV value that we calculate from the  $(i + 1)$ th frame of the trajectory will thus be similar to the value obtained for the  $i$ th trajectory frame.

9. The Monte Carlo data that was used to construct the right panel of figure 11 was generated by performing a metropolis Monte Carlo simulation that sampled points from a standard normal distribution.
10. The error bars,  $\epsilon_j$ , obtained for the components of the histogram,  $\langle g_j \rangle$ , that are calculated using equation 31 must be propagated in order to obtain the error on the free energy. As the free energy is proportional to the logarithm of the probability, the error on this quantity is  $k_B T \frac{\epsilon_j}{\langle g_j \rangle}$ .
11. When we differentiate  $\mathcal{L}$  in Eq. 35 with respect to  $p_k$  we find that at the constrained optimum:

$$\frac{\partial \mathcal{L}}{\partial p_j} = \sum_{k=1}^N \frac{t_{kj}}{p_j} + \sum_{k=1}^N \lambda_k c_k w_{kj} = 0 \quad (46)$$

which can be rearranged to give:

$$p_j = - \frac{\sum_{k=1}^N t_{kj}}{\sum_k \lambda_k c_k w_{kj}} \quad (47)$$

Similarly, when we differentiate  $\mathcal{L}$  in Eq. 35 with respect to  $c_k$  we find that

$$\frac{\partial \mathcal{L}}{\partial c_k} = \sum_{j=1}^M \frac{t_{kj}}{c_k} + \sum_{j=1}^M \lambda_k w_{kj} p_j = 0 \quad (48)$$

which rearranges to give:

$$\lambda_k = - \frac{\sum_{j=1}^M c_k w_{kj} p_j}{\sum_{j=1}^M t_{kj}} \quad (49)$$

Finally, when we differentiate  $\mathcal{L}$  in Eq. 35 with respect to  $\lambda_k$  we obtain the constraint:

$$\sum_{j=1}^M c_k w_{kj} p_j = 1 \quad (50)$$

The last two of these equations can be combined to give:

$$\lambda_k = - \frac{1}{\sum_{j=1}^M t_{kj}} \quad (51)$$

Notice that  $\sum_{j=1}^M t_{kj}$  is simply the length of trajectory  $k$ . Assuming all the trajectories to have the same length,  $\lambda_k$  will thus be independent of  $k$ . Inserting this result into equation 47 will thus give:

$$p_j \propto \frac{\sum_{k=1}^N t_{kj}}{\sum_k c_k w_{kj}} \quad (52)$$

Furthermore, rearranging equation 50 gives:

$$c_k = \frac{1}{\sum_{j=1}^M w_{kj} p_j} \quad (53)$$

## References

- [1] Massimiliano Bonomi, Davide Branduardi, Giovanni Bussi, Carlo Camilloni, Davide Provasi, Paolo Raiteri, Davide Donadio, Fabrizio Marinelli, Fabio Pietrucci, Ricardo A Broglia, et al. Plumed: A portable plugin for free-energy calculations with molecular dynamics. *Computer Physics Communications*, 180(10):1961–1972, 2009.
- [2] Gareth A Tribello, Massimiliano Bonomi, Davide Branduardi, Carlo Camilloni, and Giovanni Bussi. Plumed 2: New feathers for an old bird. *Computer Physics Communications*, 185(2):604–613, 2014.
- [3] Mark James Abraham, Teemu Murtola, Roland Schulz, Szilárd Páll, Jeremy C Smith, Berk Hess, and Erik Lindahl. Gromacs: High performance molecular simulations through multi-level parallelism from laptops to supercomputers. *SoftwareX*, 1:19–25, 2015.
- [4] Steve Plimpton. Fast parallel algorithms for short-range molecular dynamics. *Journal of computational physics*, 117(1):1–19, 1995.
- [5] Ilian T Todorov, William Smith, Kostya Trachenko, and Martin T Dove. DL\_poly\_3: new dimensions in molecular dynamics simulations via massive parallelism. *Journal of Materials Chemistry*, 16(20):1911–1918, 2006.
- [6] Jürg Hutter, Marcella Iannuzzi, Florian Schiffmann, and Joost VandeVondele. Cp2k: atomistic simulations of condensed matter systems. *Wiley Interdisciplinary Reviews: Computational Molecular Science*, 4(1):15–25, 2014.
- [7] D.A. Case, R.M. Betz, D.S. Cerutti, T.E. Cheatham, III, T.A. Darden, R.E. Duke, T.J. Giese, A.W. Goetz, H. Gohlke, N. Homeyer, S. Izadi, P. Janowski, J. Kaus, A. Kovalenko, T.S. Lee, S. LeGrand, P. Li, C. Lin, T. Luchko, R. Luo, B. Madej, D. Mermelstein, K.M. Merz, G. Monard, H. Nguyen, H.T. Nguyen, I. Omelyan, A. Onufriev, D.R. Roe, A. Roitberg, C. Sagui, C.L. Simmerling, W.M. Botello-Smith, R.C. Walker, J. Swails, J. Wang, R.M. Wolf, X. Wu, L. Xiao, and P.A. Kollman. Amber 2016, university of california, san francisco, 2016.

- [8] Peter Eastman, Jason Swails, John D Chodera, Robert T McGibbon, Yutong Zhao, Kyle A Beauchamp, Lee-Ping Wang, Andrew C Simmonett, Matthew P Harrigan, Chaya D Stern, et al. Openmm 7: rapid development of high performance algorithms for molecular dynamics. *PLoS computational biology*, 13(7):e1005659, 2017.
- [9] Giacomo Fiorin, Michael L Klein, and Jérôme Hénin. Using collective variables to drive molecular dynamics simulations. *Molecular Physics*, 111(22-23):3345–3362, 2013.
- [10] Hythem Sidky, Yamil J Colón, Julian Helfferich, Benjamin J Sikora, Cody Bezik, Weiwei Chu, Federico Giberti, Ashley Z Guo, Xikai Jiang, Joshua Lequieu, et al. Ssages: Software suite for advanced general ensemble simulations. *The Journal of Chemical Physics*, 148(4):044104, 2018.
- [11] Alejandro Gil-Ley and Giovanni Bussi. Enhanced conformational sampling using replica exchange with collective-variable tempering. *Journal of chemical theory and computation*, 11(3):1077–1085, 2015.
- [12] Robert B Best, Gerhard Hummer, and William A Eaton. Native contacts determine protein folding mechanisms in atomistic simulations. *Proceedings of the National Academy of Sciences*, 110(44):17874–17879, 2013.
- [13] Carlo Camilloni and Michele Vendruscolo. Statistical mechanics of the denatured state of a protein using replica-averaged metadynamics. *Journal of the American Chemical Society*, 136(25):8982–8991, 2014.
- [14] Yong Zhang and Gregory A Voth. Combined metadynamics and umbrella sampling method for the calculation of ion permeation free energy profiles. *Journal of chemical theory and computation*, 7(7):2277–2283, 2011.
- [15] Thierry De Meyer, Bernd Ensing, Sven MJ Rogge, Karen De Clerck, Evert Jan Meijer, and Veronique Van Speybroeck. Acidity constant (pka) calculation of large solvated dye molecules: Evaluation of two advanced molecular dynamics methods. *ChemPhysChem*, 17(21):3447–3459, 2016.
- [16] Bingqing Cheng, Gareth A Tribello, and Michele Ceriotti. Solid-liquid interfacial free energy out of equilibrium. *Physical Review B*, 92(18):180102, 2015.
- [17] Gareth A. Tribello, Federico Giberti, Gabriele C. Sosso, Matteo Salvalaglio, and Michele Parrinello. Analyzing and driving cluster formation in atomistic simulations. *Journal of Chemical Theory and Computation*, 13(3):1317–1327, 2017. PMID: 28121147.
- [18] Baron Peters. Inertial likelihood maximization for reaction coordinates with high transmission coefficients. *Chemical Physics Letters*, 554:248–253, 2012.



- [19] Wolfgang Kabsch. A solution for the best rotation to relate two sets of vectors. *Acta Crystallographica Section A: Crystal Physics, Diffraction, Theoretical and General Crystallography*, 32(5):922–923, 1976.
- [20] J. Vymetal and J. Vondrasek. Gyration-and inertia-tensor-based collective coordinates for metadynamics. application on the conformational behavior of polyalanine peptides and trp-cage folding. *The Journal of Physical Chemistry A*, 115(41):11455–11465, 2011.
- [21] Richard A Cunha and Giovanni Bussi. Unraveling mg2+-rna binding with atomistic molecular dynamics. *RNA*, 23(5):628–638, 2017.
- [22] Fabio Pietrucci and Alessandro Laio. A collective variable for the efficient exploration of protein beta-sheet structures: Application to sh3 and gb1. *Journal of Chemical Theory and Computation*, 5(9):2197–2201, 2009.
- [23] Christian Bartels and Martin Karplus. Probability distributions for complex systems: adaptive umbrella sampling of the potential energy. *The Journal of Physical Chemistry B*, 102(5):865–880, 1998.
- [24] Massimiliano Bonomi and Michele Parrinello. Enhanced sampling in the well-tempered ensemble. *Physical review letters*, 104(19):190601, 2010.
- [25] Themis Lazaridis and Martin Karplus. Effective energy function for proteins in solution. *Proteins: Structure, Function, and Bioinformatics*, 35(2):133–152, 1999.
- [26] Trang N Do, Paolo Carloni, Gabriele Varani, and Giovanni Bussi. Rna/peptide binding driven by electrostatics – insight from bidirectional pulling simulations. *Journal of chemical theory and computation*, 9(3):1720–1730, 2013.
- [27] Marco Nava, Ferruccio Palazzesi, Claudio Perego, and Michele Parrinello. Dimer metadynamics. *Journal of chemical theory and computation*, 13(2):425–430, 2017.
- [28] Sandro Bottaro, Pavel Banas, Jiri Sponer, and Giovanni Bussi. Free energy landscape of gaga and uucg rna tetraloops. *The journal of physical chemistry letters*, 7(20):4032–4038, 2016.
- [29] Vojtěch Spiwok, Petra Lipovová, and Blanka Králová. Metadynamics in essential coordinates: free energy simulation of conformational changes. *The Journal of Physical Chemistry B*, 111(12):3073–3076, 2007.
- [30] Ludovico Sutto, Marco D’Abramo, and Francesco Luigi Gervasio. Comparing the efficiency of biased and unbiased molecular dynamics in reconstructing the free energy landscape of met-enkephalin. *Journal of Chemical Theory and Computation*, 6(12):3640–3646, 2010.

- [31] Davide Branduardi, Francesco Luigi Gervasio, and Michele Parrinello. From a to b in free energy space. *The Journal of chemical physics*, 126(5):054103, 2007.
- [32] Grisell Díaz Leines and Bernd Ensing. Path finding on high-dimensional free energy landscapes. *Physical review letters*, 109(2):020601, 2012.
- [33] Vojtěch Spiwok and Blanka Králová. Metadynamics in the conformational space nonlinearly dimensionally reduced by isomap. *The Journal of chemical physics*, 135(22):224504, 2011.
- [34] D t Cremer and JA Pople. General definition of ring puckering coordinates. *Journal of the American Chemical Society*, 97(6):1354–1358, 1975.
- [35] Ming Huang, Timothy J Giese, Tai-Sung Lee, and Darrin M York. Improvement of dna and rna sugar pucker profiles from semiempirical quantum methods. *Journal of chemical theory and computation*, 10(4):1538–1545, 2014.
- [36] Massimiliano Bonomi and Carlo Camilloni. Integrative structural and dynamical biology with plumed-isdb. *Bioinformatics*, 33(24):3999–4000, 2017.
- [37] Toni Giorgino. Plumed-gui: An environment for the interactive development of molecular dynamics analysis and biasing scripts. *Computer Physics Communications*, 185(3):1109–1114, 2014.
- [38] William Humphrey, Andrew Dalke, and Klaus Schulten. Vmd: visual molecular dynamics. *Journal of molecular graphics*, 14(1):33–38, 1996.
- [39] Glenn M Torrie and John P Valleau. Nonphysical sampling distributions in monte carlo free-energy estimation: Umbrella sampling. *Journal of Computational Physics*, 23(2):187–199, 1977.
- [40] Shankar Kumar, John M Rosenberg, Djamal Bouzida, Robert H Swendsen, and Peter A Kollman. The weighted histogram analysis method for free-energy calculations on biomolecules. i. the method. *Journal of computational chemistry*, 13(8):1011–1021, 1992.
- [41] Barry Isralewitz, Sergei Izrailev, and Klaus Schulten. Binding pathway of retinal to bacterio-opsin: a prediction by molecular dynamics simulations. *Biophysical journal*, 73(6):2972–2979, 1997.
- [42] Alessandro Laio and Michele Parrinello. Escaping free-energy minima. *Proceedings of the National Academy of Sciences*, 99(20):12562–12566, 2002.
- [43] Marcella Iannuzzi, Alessandro Laio, and Michele Parrinello. Efficient exploration of reactive potential energy surfaces using car-parrinello molecular dynamics. *Physical review letters*, 90(23):238302, 2003.

- [44] Paolo Raiteri, Alessandro Laio, Francesco Luigi Gervasio, Cristian Micheletti, and Michele Parrinello. Efficient reconstruction of complex free energy landscapes by multiple walkers metadynamics. *The journal of physical chemistry B*, 110(8):3533–3539, 2006.
- [45] Giovanni Bussi, Francesco Luigi Gervasio, Alessandro Laio, and Michele Parrinello. Free-energy landscape for  $\beta$  hairpin folding from combined parallel tempering and metadynamics. *Journal of the American Chemical Society*, 128(41):13435–13441, 2006.
- [46] Stefano Piana and Alessandro Laio. A bias-exchange approach to protein folding. *The Journal of Physical Chemistry B*, 111(17):4553–4559, 2007.
- [47] Alessandro Barducci, Giovanni Bussi, and Michele Parrinello. Well-tempered metadynamics: a smoothly converging and tunable free-energy method. *Physical review letters*, 100(2):020603, 2008.
- [48] Davide Branduardi, Giovanni Bussi, and Michele Parrinello. Metadynamics with adaptive gaussians. *Journal of chemical theory and computation*, 8(7):2247–2254, 2012.
- [49] James F Dama, Michele Parrinello, and Gregory A Voth. Well-tempered metadynamics converges asymptotically. *Physical review letters*, 112(24):240602, 2014.
- [50] James F Dama, Grant Rotskoff, Michele Parrinello, and Gregory A Voth. Transition-tempered metadynamics: robust, convergent metadynamics via on-the-fly transition barrier estimation. *Journal of chemical theory and computation*, 10(9):3626–3633, 2014.
- [51] Jim Pfendtner and Massimiliano Bonomi. Efficient sampling of high-dimensional free-energy landscapes with parallel bias metadynamics. *Journal of chemical theory and computation*, 11(11):5062–5067, 2015.
- [52] Petr Hosek, Daniela Toulcova, Andrea Bortolato, and Vojtech Spiwok. Altruistic metadynamics: Multisystem biased simulation. *The Journal of Physical Chemistry B*, 120(9):2209–2215, 2016.
- [53] Fahimeh Baftizadeh, Pilar Cossio, Fabio Pietrucci, and Alessandro Laio. Protein folding and ligand-enzyme binding from bias-exchange metadynamics simulations. *Current Physical Chemistry*, 2(1):79–91, 2012.
- [54] Pratyush Tiwary and Michele Parrinello. From metadynamics to dynamics. *Physical review letters*, 111(23):230602, 2013.
- [55] Luca Maragliano and Eric Vanden-Eijnden. A temperature accelerated method for sampling free energy and determining reaction pathways in rare events simulations. *Chem Phys Lett*, 426:168–175, Jan 2006.

- [56] Jerry B. Abrams and Mark E. Tuckerman. Efficient and Direct Generation of Multidimensional Free Energy Surfaces via Adiabatic Dynamics without Coordinate Transformations. *J. Phys. Chem. B*, 112(49):15742–15757, DEC 11 2008.
- [57] Tony Lelièvre, Mathias Rousset, and Gabriel Stoltz. Computation of free energy profiles with parallel adaptive dynamics. *The Journal of Chemical Physics*, 126(13):134111, apr 2007.
- [58] Lianqing Zheng and Wei Yang. Practically efficient and robust free energy calculations: Double-integration orthogonal space tempering. *Journal of Chemical Theory and Computation*, 8(3):810–823, mar 2012.
- [59] Haohao Fu, Xueguang Shao, Christophe Chipot, and Wensheng Cai. Extended adaptive biasing force algorithm. an on-the-fly implementation for accurate free-energy calculations. *Journal of Chemical Theory and Computation*, 12(8):3506–3513, aug 2016.
- [60] Omar Valsson and Michele Parrinello. Variational approach to enhanced sampling and free energy calculations. *Phys. Rev. Lett.*, 113(9):090601, 2014.
- [61] Omar Valsson and Michele Parrinello. Well-Tempered Variational Approach to Enhanced Sampling. *J. Chem. Theory Comput.*, 11(5):1996–2002, 2015.
- [62] Andrew D White and Gregory A Voth. An Efficient and Minimal Method to Bias Molecular Simulations with Experimental Data. *Journal of Chemical Theory and Computation*, 10:3023–3030, 2014.
- [63] Glen M. Hocky, Thomas Dannenhoffer-Lafage, and Gregory A. Voth. Coarse-grained directed simulation. *Journal of Chemical Theory and Computation*, 13(9):4593–4603, 2017.
- [64] Alejandro Gil-Ley Andrea Cesari and Giovanni Bussi. Combining simulations and solution experiments as a paradigm for RNA force field refinement. *J Chem Theory Comput*, 12(12):6192–6200, dec 2016.
- [65] Andrew D White, James F Dama, and Gregory A Voth. Designing free energy surfaces that match experimental data with metadynamics. *J. Chem. Theory Comput.*, 11(6):2451–2460, 2015.
- [66] Fabrizio Marinelli and José D Faraldo-Gómez. Ensemble-biased metadynamics: A molecular simulation method to sample experimental distributions. *Biophys. J.*, 108(12):2779–2782, 2015.
- [67] Alejandro Gil-Ley, Sandro Bottaro, and Giovanni Bussi. Empirical corrections to the amber rna force field with target metadynamics. *Journal of chemical theory and computation*, 12(6):2790–2798, 2016.

- [68] Massimiliano Bonomi, Carlo Camilloni, Andrea Cavalli, and Michele Vendruscolo. Metainference: A Bayesian inference method for heterogeneous systems. *Science Advances*, 2(1):e1501177, 2016.
- [69] Christopher Jarzynski. Nonequilibrium equality for free energy differences. *Physical Review Letters*, 78(14):2690, 1997.
- [70] Massimiliano Bonomi, Alessandro Barducci, and Michele Parrinello. Reconstructing the equilibrium boltzmann distribution from well-tempered metadynamics. *Journal of computational chemistry*, 30(11):1615–1621, 2009.
- [71] Pratyush Tiwary and Michele Parrinello. A time-independent free energy estimator for metadynamics. *The Journal of Physical Chemistry B*, 119(3):736–742, 2014.
- [72] H Flyvbjerg and HG Petersen. Error estimates on averages of correlated data. *The Journal of Chemical Physics*, 91(1):461–466, 1989.
- [73] Yuji Sugita and Yuko Okamoto. Replica-exchange molecular dynamics method for protein folding. *Chemical physics letters*, 314(1):141–151, 1999.
- [74] Katsumi Murata, Yuji Sugita, and Yuko Okamoto. Free energy calculations for dna base stacking by replica-exchange umbrella sampling. *Chemical physics letters*, 385(1):1–7, 2004.
- [75] Jeremy Curuksu and Martin Zacharias. Enhanced conformational sampling of nucleic acids by a new hamiltonian replica exchange molecular dynamics approach. *The Journal of chemical physics*, 130(10):03B610, 2009.
- [76] Marc Souaille and Benoit Roux. Extension to the weighted histogram analysis method: combining umbrella sampling with free energy calculations. *Computer physics communications*, 135(1):40–57, 2001.
- [77] Michael R Shirts and John D Chodera. Statistically optimal analysis of samples from multiple equilibrium states. *The Journal of chemical physics*, 129(12):124105, 2008.
- [78] Zhiqiang Tan, Emilio Gallicchio, Mauro Lapelosa, and Ronald M Levy. Theory of binless multi-state free energy estimation with applications to protein-ligand binding. *The Journal of chemical physics*, 136(14):04B608, 2012.
- [79] Vojtěch Mlýnský, Giovanni Bussi, et al. Molecular dynamics simulations reveal an interplay between shape reagent binding and rna flexibility. *THE JOURNAL OF PHYSICAL CHEMISTRY LETTERS*, 9:313–318, 2018.
- [80] Piero Gasparotto and Michele Ceriotti. Recognizing molecular patterns by machine learning: An agnostic structural definition of the hydrogen bond. *The Journal of Chemical Physics*, 141(17):174110, 2014.

- [81] Gareth A. Tribello, Michele Ceriotti, and Michele Parrinello. Using sketch-map coordinates to analyze and bias molecular dynamics simulations. *Proceedings of the National Academy of Sciences*, 109(14):5196–5201, 2012.
- [82] Gareth A. Tribello, Michele Ceriotti, and Michele Parrinello. Using sketch-map coordinates to analyze and bias molecular dynamics simulations. *Proceedings of the National Academy of Sciences*, 109(14):5196–5201, 2012.
- [83] Mohammad M. Sultan and Vijay S. Pande. tica-metadynamics: Accelerating metadynamics by using kinetically selected collective variables. *Journal of Chemical Theory and Computation*, 13(6):2440–2447, 2017. PMID: 28383914.
- [84] W. Chen and A. L Ferguson. Molecular enhanced sampling with autoencoders: On-the-fly collective variable discovery and accelerated free energy landscape exploration. *ArXiv e-prints*, December 2018.
- [85] M. M. Sultan, H. K. Wayment-Steele, and V. S. Pande. Transferable neural networks for enhanced sampling of protein dynamics. *ArXiv e-prints*, January 2018.
- [86] Open path sampling. <http://openpathsampling.org/latest>.
- [87] MBBJM Tuckerman, Bruce J Berne, and Glenn J Martyna. Reversible multiple time scale molecular dynamics. *The Journal of chemical physics*, 97(3):1990–2001, 1992.
- [88] Marco Jacopo Ferrarotti, Sandro Bottaro, Andrea Pérez-Villa, and Giovanni Bussi. Accurate multiple time step in biased molecular simulations. *Journal of chemical theory and computation*, 11(1):139–146, 2014.