

## Consistent principal component modes from molecular dynamics simulations of proteins

Rodrigo Cossio-Pérez, Juliana Palma, and Gustavo Pierdominici-Sottile

*J. Chem. Inf. Model.*, **Just Accepted Manuscript** • DOI: 10.1021/acs.jcim.6b00646 • Publication Date (Web): 16 Mar 2017

Downloaded from <http://pubs.acs.org> on March 17, 2017

### Just Accepted

“Just Accepted” manuscripts have been peer-reviewed and accepted for publication. They are posted online prior to technical editing, formatting for publication and author proofing. The American Chemical Society provides “Just Accepted” as a free service to the research community to expedite the dissemination of scientific material as soon as possible after acceptance. “Just Accepted” manuscripts appear in full in PDF format accompanied by an HTML abstract. “Just Accepted” manuscripts have been fully peer reviewed, but should not be considered the official version of record. They are accessible to all readers and citable by the Digital Object Identifier (DOI®). “Just Accepted” is an optional service offered to authors. Therefore, the “Just Accepted” Web site may not include all articles that will be published in the journal. After a manuscript is technically edited and formatted, it will be removed from the “Just Accepted” Web site and published as an ASAP article. Note that technical editing may introduce minor changes to the manuscript text and/or graphics which could affect content, and all legal disclaimers and ethical guidelines that apply to the journal pertain. ACS cannot be held responsible for errors or consequences arising from the use of information contained in these “Just Accepted” manuscripts.

# CONSISTENT PRINCIPAL COMPONENT MODES FROM MOLECULAR DYNAMICS SIMULATIONS OF PROTEINS

Rodrigo Cossio-Pérez, Juliana Palma,<sup>\*</sup> and Gustavo Pierdominici-Sottile

*Departamento de Ciencia y Tecnología, Universidad Nacional de Quilmes, Sáenz Peña 352,  
B1876BXD Bernal, Argentina.*

E-mail: juliana@unq.edu.ar

## Abstract

Principal component analysis is a technique widely used for studying the movements of proteins using data collected from molecular dynamics simulations. In spite of its extensive use the technique has a serious drawback: equivalent simulations do not afford the same PC-modes. In this article we show that concatenating equivalent trajectories and calculating the PC-modes from the concatenated one significantly enhances the reproducibility of the results. Moreover, the consistency of the modes can be systematically improved by adding more individual trajectories to the concatenated one.

---

<sup>\*</sup>To whom correspondence should be addressed

## Introduction

The understanding of how proteins work needs a joint description of their dynamical and structural characteristics. Molecular dynamics (MD) simulations constitute a powerful approach to investigate their dynamical features. Within this framework, the use of principal component analysis (PCA) has emerged as one of the most widely employed techniques to analyze protein movements.<sup>1-4</sup> The methodology was introduced in MD studies of proteins by Karplus *et. al.*,<sup>5,6</sup> but it has also been extensively used in other branches of science and technology.<sup>7,8</sup>

In MD studies, PCA is mainly applied to change from a description based on local atomic coordinates to one provided by collective coordinates, called the PC-modes. They describe the simultaneous motion of different parts of the protein. The PC-modes are the eigenvectors of the covariance matrix, which is calculated with the configurations sampled from a MD trajectory.<sup>1</sup> A key feature of PCA is that only a bunch of PC-modes, those associated with the highest eigenvalues, account for approximately 70% to 90% of the total fluctuations of the protein.<sup>9-11</sup> This allows for a huge reduction in the number of degrees of freedom required to indicate the deformations of the system: just a few PC-modes provide a description equivalent to hundreds or thousands of atomic coordinates. The subspace formed by the PC-modes associated to the largest eigenvalues is called the essential space (ES). It is usually assumed that the motions related to the biological function of a protein are contained within its ES.<sup>12-15</sup> The remaining PC-modes account for irrelevant, small-amplitude fluctuations. It is said that they span the “near-constraint subspace” which is normally of no interest.<sup>16</sup>

In spite of the many examples in which PCA has proved to be useful, the PC-modes calculated by the standard procedure have an undesirable characteristic that casts doubts on their actual significance and utility: equivalent simulations do not afford the same PC-modes.<sup>17-19</sup> This lack of reproducibility can be assessed by calculating the inner product between PC-modes obtained by equivalent but independent MD simulations. Ideally the

1  
2  
3 absolute value of these inner products should be 1.0 for modes having the same index and  
4 zero otherwise. A less restrictive condition consists of requiring that the ESs obtained by  
5 equivalent simulations describe the same subspace. The parameter that evaluates the overlap  
6 between such subspaces is called the root mean squared inner product (RMSIP), and has  
7 been widely used to assess the convergence of the main PC-modes of proteins. The RMSIP  
8 is equal to 1.0 if the two ESs span the same subspace while it is zero if they are orthogonal.  
9 It is also a common practice to calculate the RMSIP with the ESs obtained from different  
10 halves of the same trajectory, using increasingly longer simulation times. This is done to  
11 evaluate the convergence of the main PC-modes with respect to the length of the simulation.  
12 This test has been applied to many different systems and in all cases the same qualitative  
13 result was obtained.<sup>20-26</sup> Initially the RMSIP grows fast, but then it levels out reaching a  
14 plateau value which is sensibly smaller than 1.0. This behavior indicates that extending the  
15 simulation time, within the ranges typically used in current MD simulations, is not effective  
16 to improve the convergence of the ES.

17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32 Many studies have focused on evaluating the main features of protein PC-modes.<sup>27-31</sup>  
33 They comprise from small soluble proteins to large membrane proteins, and from short  
34 simulations of just a few ns to relatively long simulations of more than 50 ns. In all cases it  
35 was found that the main PC-modes obtained by equivalent trajectories were different, and  
36 that the RMSIP computed from them was smaller than one. In this article we show that the  
37 consistency of the PC-modes can be improved by employing a correlation matrix obtained by  
38 concatenating independent but equivalent trajectories. The performance of the procedure  
39 is demonstrated by applying it to the principal component analysis of bovine pancreatic  
40 trypsin inhibitor (BPTI) and lysozyme.

41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51 The use of PCA of concatenated trajectories was introduced by Berendsen and co-workers  
52 in 1995.<sup>32</sup> Since then, it has been widely used as one of the diverse tools employed to charac-  
53 terize the most relevant motions of proteins and other systems of biological interest. However,  
54 the results of those studies were interpreted on intuitive foundations since analytical formulas  
55  
56  
57  
58  
59  
60

1  
2  
3 showing the precise meaning of the eigenvectors and eigenvalues obtained by concatenated-  
4 PCA had not been provided. Recently we presented such analytical expressions and discussed  
5 that two opposite limits could be found.<sup>24</sup> One extreme case occurs when the trajectories  
6 belong to two or more free energy minima, and the root mean square deviations between  
7 these minima are significantly larger than the typical fluctuations around them (i.e. the tra-  
8 jectories cannot go from one minimum to the other). This case was thoroughly analyzed  
9 in Ref. 24. The other extreme occurs when the concatenated trajectories are initiated and  
10 remain within the same free energy well. The method proposed in this article is aimed to  
11 univocally determine the main PC-modes in this case. We note that only in this situation the  
12 ES of a protein contains a set of collective coordinates useful to describe the most important  
13 fluctuations. When multiple minima are implied the main PC-modes contain the so-called  
14 “static contribution”.<sup>24</sup> It can partially or completely mask the “dynamic contribution” of  
15 protein fluctuations.  
16  
17

18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30 Concatenated trajectories have also been used for other purposes. In particular, they were  
31 employed in different methodologies developed to characterize protein folding processes. In  
32 these methods, the concatenation has to be performed following specific prescriptions aimed  
33 to accelerate the conformational sampling. Important examples of these applications can be  
34 found in Ref. 33 - 34 and the references cited therein. The analysis of the PC-modes obtained  
35 with this kind of trajectories lies outside the scope of this article. Instead we concentrate  
36 on the characterization of the dynamics of stable conformations of proteins, whose main  
37 PC-modes have a clear and straightforward interpretation as collective coordinates useful to  
38 describe their most important fluctuations.  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50

## 51 **Materials and Methods**

### 52 **Statement of the problem**

53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 In applications of PCA to the study of protein dynamics, the PC-modes are obtained by  
4 diagonalizing a correlation matrix whose elements are given by,  
5  
6  
7

$$C_{ij} = \frac{1}{N} \sum_{k=1}^N (x_i^k - \bar{x}_i)(x_j^k - \bar{x}_j). \quad (1)$$

8  
9  
10 Here  $x_i^k, x_j^k$ , are a pair of elements of vector  $\mathbf{x}^k$ , which describes the configuration of the  
11 system at time step  $k$ , while  $\bar{x}_i$  and  $\bar{x}_j$  are their average values calculated from the  $N$   
12 structures sampled in the MD simulation. Normally,  $\mathbf{x}$  is a vector containing the Cartesian  
13 coordinates of the  $C_\alpha$  atoms of the protein, but other choices can also be used.<sup>35-37</sup> During  
14 the setting of any MD simulation several parameters are chosen at random. Besides, the  
15 sampling interval and the time at which the first structure is recorded are arbitrarily decided.  
16 Therefore the atomic coordinates of the selected structures are random numbers too, as are  
17 random the  $C_{ij}$  coefficients calculated from them.  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28

29 For an infinite long simulation, from which an infinite number of samples could be taken,  
30 the  $C_{ij}$  coefficients would assume perfectly defined values, free of statistical errors. However,  
31 infinite long simulations are not possible and the experience indicates that currently feasi-  
32 ble simulations are not long enough to produce correlation matrices with sufficiently small  
33 statistical uncertainties. Thus, if  $\mathbf{C}^\infty$  is the correlation matrix corresponding to an infinite  
34 long sampling and  $\mathbf{C}$  is a correlation matrix computed from a finite number of samples, the  
35 relation between them can be written as,  
36  
37  
38  
39  
40  
41  
42  
43  
44

$$\mathbf{C} = \mathbf{C}^\infty + \mathbf{E}, \quad (2)$$

45  
46  
47  
48 where the elements of matrix  $\mathbf{E}$  contain the statistical errors of the correlation coefficients  
49 computed from the sample. Since the elements of  $\mathbf{E}$  are different than zero, the eigenvalues  
50 and eigenvectors of matrix  $\mathbf{C}$  differ from those of the actual correlation matrix  $\mathbf{C}^\infty$ . These  
51 discrepancies can be estimated if the elements of  $\mathbf{E}$  are small enough and fulfill some defined  
52 characteristics (i.e. they are normally distributed).<sup>38</sup> Such estimations have successfully been  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 done in applications of PCA to other branches of science, technology and economy.<sup>39,40</sup>

4  
5 Recognizing that infinite sampling is not possible, one faces the problem of determining  
6 which sampling procedure minimizes the statistical error of the results. The use of increas-  
7 ingly longer simulations is normally employed but, as already described in the introduction,  
8 such strategy does not lead to the desired result. Alternatively, one can use several equiv-  
9 alent MD simulations that just differ from each other in the initial velocities of the atoms.  
10 We recently demonstrated that the correlation matrix  $\mathbf{C}^{(n)}$ , obtained by concatenating  $n$   
11 trajectories with the same number of samples can be decomposed as,<sup>24</sup>

$$\mathbf{C}^{(n)} = \frac{1}{n} \sum_{i=1}^n \mathbf{C}_i + \mathbf{S}^{(n)}. \quad (3)$$

12 Here  $\mathbf{C}_i$  is the correlation matrix corresponding to the  $i$ -th trajectory while  $\mathbf{S}^{(n)}$  is the  
13 correlation matrix computed from the  $n$  average structures. If the individual MD simulations  
14 were able to sample all regions of the accessible configurational space, trajectories that just  
15 differ in their initial atomic velocities would produce almost the same average structures.  
16 In this case matrix  $\mathbf{S}^{(n)}$  should be significantly smaller than the  $\mathbf{C}_i$  matrices, because the  
17 deviation of the individual average structures with respect to the global average would be  
18 much smaller than the fluctuations observed in any single trajectory. Under such conditions,  
19 the correlation matrix of the concatenated trajectory becomes quite close to the average of  
20 the individual correlation matrices,  $\mathbf{C}_{av}^{(n)}$ .

$$\mathbf{C}^{(n)} \approx \frac{1}{n} \sum_{i=1}^n \mathbf{C}_i = \mathbf{C}_{av}^{(n)}. \quad (4)$$

21 According to the Classical Central Limit Theorem the last equality implies that for large  
22 values of  $n$ , the statistical uncertainty in the elements of  $\mathbf{C}_{av}^{(n)}$  will be smaller than those  
23 of the  $\mathbf{C}_i$ 's by a factor of about  $1/\sqrt{n}$ . Thus, one is induced to think that concatenating  
24 trajectories should provide a route for obtaining reproducible PC-modes.

25 However, it could also happen that the single trajectories are relatively short or get

1  
2  
3  
4 trapped and sample only a portion of the available configurational space. In such case,  
5  
6 the individual averages will be disperse with respect to the global average structure and  
7  
8 matrix  $\mathbf{S}^{(n)}$  will not be negligible. Even in this apparently adverse condition, the correlation  
9  
10 matrix of the concatenated trajectory can still be equal to the average of some conveniently-  
11  
12 defined correlation matrices  $\mathbf{C}_i$ . To see how to achieve this, one has to recognize that matrix  
13  
14  $\mathbf{C}^{(n)}$  remains unchanged with respect to any permutation of the structures used to compute  
15  
16 it. Usually, structures 1 to  $N$  correspond to the first simulation, structures  $N + 1$  to  $2N$   
17  
18 to the second one, and so on. However one can shuffle the  $nN$  structures employed to  
19  
20 calculate  $\mathbf{C}^{(n)}$  and then divide them into  $n$  sets of  $N$  arbitrarily-selected structures. Any  
21  
22 of these new sets will have structures originated from different MD simulations. Therefore,  
23  
24 for sufficiently large  $n$  and  $N$ , the average structures of these sets will be pretty similar to  
25  
26 each other, and the new  $\mathbf{S}^{(n)}$  matrix will be negligible with respect to the new  $\mathbf{C}_{av}^{(n)}$ . Thus,  
27  
28 the shuffling procedure does not affect  $\mathbf{C}^{(n)}$  but it changes  $\mathbf{C}_{av}^{(n)}$  and  $\mathbf{S}^{(n)}$ , in such a way  
29  
30 that their changes mutually compensate. When the structures sampled from a single MD  
31  
32 simulation are biased, the correlation matrix computed from them underestimates the actual  
33  
34 correlations. This occurs because the deviations of the sampled structures with respect to  
35  
36 their own average are smaller than their deviations with respect to the true average. For  
37  
38 the same reason, also matrix  $\mathbf{C}_{av}^{(n)}$  underestimates the actual correlations. The calculation of  
39  
40 matrix  $\mathbf{C}^{(n)}$  corrects this error because the correlations that get lost in  $\mathbf{C}_{av}^{(n)}$  appear in  $\mathbf{S}^{(n)}$ .  
41  
42 We therefore conclude that, even if the individual trajectories perform a biased sampling of  
43  
44 the available configurational space, matrix  $\mathbf{C}^{(n)}$  converges to an average of  $n$  conveniently-  
45  
46 defined individual correlation matrices. Accordingly, its statistical uncertainty is reduced by  
47  
48 a factor of  $1/\sqrt{n}$ .  
49

50  
51 In what follows, we will refer to matrix  $\mathbf{C}^{(n)}$  as the correlation matrix of the concatenated  
52  
53 trajectory. However, from the previous discussion, it should be clear that there is not a real  
54  
55 need to “concatenate” the trajectories, since any order of the whole set of structures produces  
56  
57 the same result. The key point here is to compute the correlation matrix using structures  
58  
59  
60



1  
2  
3 sampled from multiple equivalent trajectories.  
4  
5  
6

## 7 **Molecular dynamics simulations**

8

9  
10 Principal component analysis of BPTI and lysozyme were used to test the consistency of the  
11 PC-modes. The initial coordinates of the proteins were obtained from the Protein Data Bank,  
12 ID=5PTI for BPTI<sup>41</sup> and 1REX for lysozyme.<sup>42</sup> The systems were solvated in a truncated  
13 octahedral cell of TIP3P explicit water molecules and minimized at constant volume. In a  
14 second stage they were heated at constant volume from 0 K to 310 K during 1 ns, using the  
15 weak coupling algorithm with  $\tau_P=2.8$  ps. After that, we switched to constant temperature  
16 and pressure conditions using a value of 2.0 ps for both,  $\tau_{TP}$  and  $\tau_P$ . Finally, an equilibration  
17 run of 10 ns was performed. For each system, the final structure of the equilibration stage was  
18 used as the initial configuration of the production runs. We run 180 equivalent trajectories  
19 of 5 ns and 80 trajectories of 50 ns for each system. These trajectories just differed in the  
20 initial velocities of the atoms, which were chosen from a maxwellian distribution at 310 K.  
21 Snapshots were taken every 25 ps in the 5-ns trajectories and every 250 ps in the 50-ns  
22 trajectories.  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35

36 The simulations were performed with the AMBER 14 package using the ff99SB force  
37 field, applying periodic boundary conditions with a cutoff of 12.0 Å. The shake algorithm  
38 was employed to maintain bond distances to hydrogen, allowing for a time step of 2.0 fs.  
39 We tested that the projections of the individual trajectories onto their first two PC-modes  
40 did not resemble cosine functions.<sup>43</sup> Also, it was checked that the RMSIP calculated from  
41 different halves of the same trajectory was converged with respect to the simulation time.<sup>20</sup>  
42 Thus, any of the individual simulations employed in this work would pass the convergence  
43 assessments usually applied in PCA studies of protein.  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## Measurements of the convergence

The convergence of PCA is usually evaluated by comparing the PC-modes obtained from a reference trajectory against the ones derived from trial trajectories whose convergence one is trying to evaluate. Normally the reference trajectory is much longer than the trial trajectories or it has been obtained with procedures that improve the conformational sampling. This kind of assessment is based on the assumption that the reference trajectory provides a fair enough sampling, so that its PC-modes are already converged. Here we apply a criterion that does not require the *a-priori* knowledge of such reference trajectory. Instead, we consider that the PC-modes or the essential space of a protein are converged when two alternative but otherwise equivalent computations afford the same PC-modes or ESs. The similarity between the PC-modes and ESs of the alternative computations were measured using different parameters. They are described in the following paragraphs.

We analyzed the absolute value of the scalar product between corresponding PC-modes calculated from alternative simulations *a* and *b*,  $|\text{PC}_i^a \cdot \text{PC}_i^b|$ , with particular emphasis on the first mode,  $\text{PC}_1$ . We also employed the RMSIP, which measures the common portion of the ESs determined from the pair MD simulations,

$$\text{RMSIP}_M = \sqrt{\frac{1}{M} \sum_{i=1}^M \sum_{j=1}^M |\text{PC}_i^a \cdot \text{PC}_j^b|^2}. \quad (5)$$

Here  $M$  denotes the dimension of the subspaces while  $\text{PC}_i^a$  and  $\text{PC}_j^b$  are the  $i$ -th and  $j$ -th eigenvectors obtained from simulations *a* and *b*, respectively. In our analysis, we set  $M = 2$  since it is a common practice to analyze protein motions in the subspace spanned by the first two eigenvectors.<sup>18,25,28,44–46</sup> It should be noted that the use of  $M = 2$  makes the evaluation as strict as possible. Finally, for each pair of trajectories, we also evaluated the covariance overlap,  $s$ , proposed by Hess.<sup>43</sup> The overlap is not a measure of the convergence of the essential space. Instead, it assesses the similarity of the spaces sampled from a given pair of

trajectories. The overlap is defined as,

$$s = 1 - d_N(\mathbf{C}_a, \mathbf{C}_b), \quad (6)$$

where  $d_N$  is the normalized distance between the correlation matrices  $\mathbf{C}_a$  and  $\mathbf{C}_b$ . The normalized distance is calculated as

$$d_N(\mathbf{C}_a, \mathbf{C}_b) = \left[ \frac{\text{tr}(\mathbf{C}_a + \mathbf{C}_b - 2\mathbf{C}_a^{1/2}\mathbf{C}_b^{1/2})}{\text{tr}(\mathbf{C}_a) + \text{tr}(\mathbf{C}_b)} \right]^{1/2}, \quad (7)$$

where  $\mathbf{C}_\alpha^{1/2}$ , the square root of correlation matrix  $\mathbf{C}_\alpha$ , is calculated as,

$$\mathbf{C}_\alpha^{1/2} = \mathbf{R}_\alpha \mathbf{\Lambda}_\alpha^{1/2} \mathbf{R}_\alpha^T. \quad (8)$$

In the last equation  $\mathbf{R}_\alpha$  is the matrix that diagonalizes  $\mathbf{C}_\alpha$  and  $\mathbf{\Lambda}_\alpha$  is the diagonal matrix that contains its eigenvalues. When two independent simulations afford the same sampling, the two correlation matrices are the same, the distance between them is zero and the overlap amounts to 1.0. On the contrary, if the subspaces sampled in the simulations are orthogonal, the normalized distance evaluates to 1.0 and the overlap is zero.

## Results

Each individual simulation was used to compute a set of PC-modes. Therefore, for each system, we obtained 180 sets of PC-modes with the 5-ns trajectories and 80 sets with the 50-ns trajectories. For both, BPTI and lysozyme, PC-mode sets computed from trajectories of equal length were grouped into pairs. All possible pairs were generated. Thus, we formed 16110 pairs with the trajectories of 5 ns and 3160 pairs with the trajectories of 50 ns. For each pair we calculated the absolute value of the inner product  $|\text{PC}_i^a \cdot \text{PC}_i^b|$ , the RMSIP<sub>2</sub> and the overlap. This allowed us to estimate reliable probability distributions for the three

parameters.

To evaluate the hypothesis that concatenating trajectories improves the reproducibility of the PC-modes we collected individual simulations of equal length into batches of  $n$  simulations, so that each trajectory was allocated into a single batch. Then, we concatenated the trajectories of a given batch and computed  $\mathbf{C}^{(n)}$ , the correlation matrix of the concatenated trajectory. Finally, the reproducibility of the PC-modes so obtained was assessed by computing  $|\text{PC}_i^a \cdot \text{PC}_i^b|$ , RMSIP<sub>2</sub> and overlap, for all possible pairs of batches formed with the given  $n$ . We tried different values of  $n$ . In Table 1 we present the alternative values of  $n$  employed in this work, along with the number of batches and the number of pairs of batches that can be formed with the given  $n$ .

## PC modes from single trajectories of BPTI

Figure 1 shows the probability distributions for  $|\text{PC}_1^a \cdot \text{PC}_1^b|$ , calculated from all the independent simulations of BPTI. The vertical black line indicates the value of the inner product that corresponds to 99% of cumulative probability, for normalized random vectors of the same size. The cumulative probability was evaluated as,

$$P_{\text{cum}}(x^*) = \int_0^{x^*} \rho(x) dx, \quad (9)$$

where  $\rho(x)$  is the probability density that a random vector of dimension  $M$  has a square projection  $x$  onto a subspace of dimension  $m$ . For the present case  $M = 58 \times 3 = 174$ , since the model of BPTI has 58 residues, while  $m = 1$  since we are considering the projection onto a single PC-mode. According to the equations provided by Amadei et. al. in Ref. 20,  $\rho(x)$  is given by,

$$\rho(x) = \frac{(M-1)!}{(m-1)!(M-m-1)!} x^{(m-1)} (1-x)^{(M-m-1)}, \quad (10)$$

that in our case simplifies to  $\rho(x) = (M-1)(1-x)^{(M-2)}$ .

It is seen that the scalar product between the first PC-modes of individual MD simulations

1  
2  
3 affords significantly larger values than those expected for random vectors. However very low  
4 values, which imply almost orthogonal  $PC_1$ , are likely to be obtained too. Figure 1 also shows  
5 that the reproducibility of  $PC_1$  is better for the trajectories of 50 ns than for those of 5 ns.  
6 However, even for these longer simulations, pretty low values of the scalar product are usually  
7 found. The distributions presented in Fig. 1 remind the ones reported by Grossfield *et. al.*<sup>31</sup>  
8 who analyzed the reproducibility of  $PC_1$  for different membrane systems, using relatively  
9 long trajectories. Thus, Fig. 1 does not provide new evidence about the characteristics of  
10  $PC_1$  obtained from individual MD simulations, but just reinforces the conclusion that its  
11 direction is random to a large extent.  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21

22 Fig. 2 presents typical examples for the inner-product matrices formed with the first six  
23 PC-modes of two equivalent simulations. Ideally, the out-of-diagonal elements should be  
24 null. The pictures presented in Fig. 2 differ from this ideal, indicating that the first PC-  
25 modes of a given simulation are mostly distributed among the first PC-modes of the other  
26 simulation. This characteristic has already been described for other systems.<sup>27,47</sup> Because  
27 of this behavior, it is more difficult to obtain a one-to-one correspondence between the  
28 individual PC-modes than to converge the subspace spanned by some of them. This is in  
29 fact the conclusion afforded by the probability distributions of  $RMSIP_2$ , presented in Fig. 3,  
30 which show that pretty low values of  $RMSIP$  are rather unlikely. However, the agreement  
31 for  $RMSIP_2$  is still far from satisfactory since the most likely values are just  $\approx 0.55$  (5-ns  
32 trajectories) and  $\approx 0.65$  (50-ns trajectories). The probability distributions for the covariance  
33 overlap are also presented in Fig. 3. They are somewhat narrower than those of  $RMSIP_2$  and  
34 their maxima are shifted to the left. In spite of these differences the conclusions attained from  
35 both distributions,  $RMSIP_2$  and overlap, are similar. We finally note that the distributions  
36 obtained with the longer trajectories are shifted to the right of those computed with the  
37 shorter ones. Thus, increasing the simulation time helps to improve the reproducibility of  
38 the essential space but, as noted above, the results are still far from satisfactory.  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## PC modes from the concatenated trajectories of BPTI

Figure 4 shows the evolution of the overlap and  $\text{RMSIP}_2$  as the number of simulations in the concatenated trajectory,  $n$ , is increased. For the short trajectories one can readily appreciate that both, overlap and  $\text{RMSIP}_2$ , raise very rapidly between 1 and 30 trajectories. Then there is a change of behavior and the two parameters variate more slowly. In particular, the overlap seems to level out at about  $n = 60$  where it reaches a value of 0.88. The increase of  $\text{RMSIP}_2$ , on the other hand, does not stop after  $n = 60$  but just becomes slower. For  $n = 90$  we found a  $\text{RMSIP}_2$  of 0.98 indicating that the subspaces spanned by the first two PC-modes of such concatenated trajectories are nearly the same.  $\text{RMSIP}_2$  and overlap of the long trajectories also increase with  $n$ . However, in this case, the initial improvement is not so marked as in the case of the short trajectories. This is mainly because the two parameters start from higher average values. In this case, we found a  $\text{RMSIP}_2$  of 0.98 for  $n = 40$ . For this  $n$  the overlap evaluates to 0.89.

The fact that the  $\text{RMSIP}_2$  is almost fully converged while the overlap is not indicates that the first two PC-modes are mostly contained in the subspace shared by the two concatenated trajectories, while the orthogonal subspace accounts for the non-important PC-modes. Besides, this reveals that the first PC-modes converge faster than those corresponding to higher indexes. This is really a good characteristic since PCA is normally used to reduce the dimensionality of the system under analysis. The concatenated-PCA technique can reliably determine the subspace that contains the most important collective motions of the protein, even though the sampling of the available configurational space is still not perfect.

In general, for a given  $n$ ,  $\text{RMSIP}_2$  calculated from the 50-ns trajectories gets higher averages than those derived from the 5-ns trajectories. For example, for  $n = 5$ , the average of  $\text{RMSIP}_2$  is 0.68 for the short trajectories and 0.78 for the long ones. For  $n = 20$ , the values are 0.82 and 0.90, respectively. However, the benefits of concatenating trajectories become more evident if one considers not just the average values but the whole range of possible outcomes, for an equivalent simulation time. Thus for example, by concatenating

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

10 trajectories of 5 ns we obtained an average RMSIP<sub>2</sub> of 0.74. Exactly the same average was obtained for single trajectories of 50 ns. But, in the former case, the range of possible outcomes was 0.53-0.96 while in the second one it was 0.09-0.96. A closer inspection to the results presented in Fig. 4 shows that the increase in the average values of RMSIP<sub>2</sub> with  $n$  is mainly caused by the increase in the minimum values while the maxima are all similar. Thus, by doing a principal component analysis with concatenated trajectories one can avoid obtaining really ill-defined PC-modes.

It remains to be checked how the reproducibility of the individual PC-modes varies with the number of simulations included in the concatenated trajectory. In the following, we will just present the results obtained with the 5-ns trajectories. The previous discussion demonstrates that it is more difficult to obtain consistent PC-modes in this case than with the longer trajectories. Figure 5 shows typical inner product matrices for the first six PC-modes obtained from  $\mathbf{C}^{(n)}$ , for  $n = 10, 20, 30$  and 90. For  $n = 90$  we only have one matrix to show. For  $n = 10, 20$  and 30 we have several options and decided to show matrices that present an average behavior. This is, neither the best nor the worst matrix for the given value of  $n$ , but an intermediate one. The improvement in the reproducibility of the individual PC-modes can be clearly seen. By increasing  $n$ , the elements in the diagonal or their closest neighbors reach significant values, while more distant elements become smaller.

For  $n = 90$ , the inner product between the two PC<sub>1</sub> is 0.979, while that of the PC<sub>2</sub> is 0.977. Similar almost perfect agreement is found for PC<sub>5</sub> and PC<sub>6</sub>. However, something odd seems to happen between the PC<sub>3</sub> and PC<sub>4</sub> since the PC<sub>3</sub> of one batch is mostly contained in PC<sub>4</sub> of the other one and vice versa. This occurs because PC<sub>3</sub> and PC<sub>4</sub> are almost degenerate. The direction of degenerate eigenvectors is arbitrary since any linear combination of such vectors is also an eigenvector with the same eigenvalue. Thus, in cases like this, one cannot do better than determining the subspace spanned by the degenerate vectors. The near-constrained subspace is plenty of almost degenerate eigenvectors. However, this causes no troubles because these vectors are discarded in applications of PCA to molecular dynamics

1  
2  
3 of proteins. On the contrary, if the degeneracy occurs in vectors with high eigenvalues, all of  
4 them have to be included in the essential subspace. It is important to note that the mixing  
5 between  $PC_3$  and  $PC_4$  shown in Fig. 5 is fortuitous. We could have missed it if we had  
6 grouped the simulations in different batches of 90 elements.  
7  
8  
9  
10

## 11 **PC modes for the concatenated trajectories of lysozyme**

12  
13 The distributions of  $|PC_1^a \cdot PC_1^b|$ ,  $RMSIP_2$  and overlap, computed from individual trajectories  
14 of lysozyme follow the same qualitative behavior as those of BPTI (presented in Figs. 1 and  
15 3). They show that PC-modes calculated from a single simulation, either of 5 ns or 50 ns, are  
16 poorly defined. To be concise we will not present those distributions here. Instead we will  
17 focus on what happens when trajectories are concatenated since that is the main subject  
18 of this work. We show in Fig. 6 the evolution of  $RMSIP_2$  and overlap as the number of  
19 concatenated trajectories,  $n$ , is increased. It is readily noted that the same trends observed  
20 for BPTI also apply to lysozyme. The average values of  $RMSIP_2$  and overlap increase with  
21  $n$ . This is mainly caused by the increase of the minimum possible values while the maximum  
22 values are all similar and high. For the same  $n$ , better results are obtained with trajectories  
23 of 50 ns than those of 5 ns. However, what is more relevant here, is the comparison of  
24 the results corresponding to the same simulation time. For example, the average  $RMSIP_2$   
25 between independent trajectories of 50 ns is 0.64, with a lower bound of 0.13 and an upper  
26 bound of 0.93. On the other hand, if one calculates  $RMSIP_2$  between trajectories obtained  
27 by concatenating 10 simulations of 5 ns, the average value is 0.82 and the boundaries are  
28 0.57 and 0.97. Thus, concatenating trajectories improves the reproducibility of the results  
29 and significantly reduces the chances of getting ill-defined PC-modes. As observed in the  
30 case of BPTI, for the highest  $n$  tried in this work,  $RMSIP_2$  is very close to 1.0 while the  
31 overlap still noticeably deviates from that. This reinforces the observation that the first PC-  
32 modes converge faster than those corresponding to higher indexes. Therefore, well-defined  
33 essential spaces for proteins can be computed even though the sampling of the available  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60



1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

configurational space is not fully-converged.

## Discussion

The results presented above demonstrate that reproducible PC-modes can be obtained by diagonalizing the correlation matrix of a concatenated trajectory, formed from  $n$  equivalent but independent simulations. Agreement of the most important PC-modes, within any desired accuracy, can be reached by systematically increasing  $n$ . A particularly convenient characteristic of the procedure is that relatively small values of  $n$  (i.e. between 10 and 20) significantly increase the average RMSIP by elevating the minimum possible value. In this way, one can avoid obtaining ill-defined essential spaces. The convergence of the individual PC-modes, on the other hand, requires larger values of  $n$ .

The reason why concatenating trajectories improves the reproducibility of the PC-modes was outlined in Section “Statement of the problem” and is based on the formulas presented in Ref. 24. A practical example of the performance of the procedure can be seen in Fig. 7. It shows the projection of typical 5-ns trajectories of BPTI onto the plane spanned by the first two eigenvectors of  $\mathbf{C}^{(180)}$ . A contour plot of the free energy calculated with the whole set of 5-ns trajectories of BPTI is also shown there. It is observed that single trajectories just occupy a fraction of the available area. They can repeatedly pass through a given region and never visit a nearby accessible zone. The use of different initial velocities takes the trajectories to different regions. Therefore, even though they individually move around a limited zone, the full accessible region is recovered when they are considered altogether.

Many years ago Caves et. al. described the same behavior for simulations of crambin<sup>18</sup>. They concluded that multiple equivalent trajectories are more efficient to provide a fair sampling of the available conformational space than a single long trajectory. Other studies of that time made similar observations and attained to the same conclusion<sup>48–53</sup>. However the

1  
2  
3 use of multiple trajectories is normally not considered a requirement for obtaining statistically  
4 significant results. Therefore it is not routinely used. Very recent investigations are re-  
5 installing the subject<sup>54,55</sup>. The present work aims to contribute in the same direction, but  
6 the focus is put on the behavior of PC-modes. To the best of our knowledge, this issue has  
7 never been systematically studied before.  
8  
9

10  
11  
12  
13  
14 Finally, it is interesting to discuss how the PC-modes obtained by concatenating short  
15 trajectories compare with those calculated by concatenating long trajectories. Unfortunately,  
16 there is not a single answer to that question since it depends on several factors such as the  
17 time-scale of the simulations being considered or the rigidity of the structure under analysis.  
18 Ref. 18 described that trajectories starting from the same structure but differing in the  
19 atomic velocities rapidly diverge from the initial point. Then they stabilize and move within  
20 a hyperspherical cortex defined by a nearly constant RMSD with respect to the original  
21 structure<sup>18</sup>. If that were strictly the case, converged PC-modes computed by concatenating  
22 short and long trajectories would be the same. However more recent studies showed that the  
23 description of Ref. 18 does not hold on the much longer time-scales affordable nowadays<sup>56</sup>.  
24 In general, if during the extra time the trajectories move further away from their initial point,  
25 converged PC-modes computed by concatenating short and long trajectories will differ. In  
26 the present study, we found that RMSIP<sub>2</sub> for converged PC-modes computed from 5-ns  
27 and 50-ns trajectories was 0.595 for BPTI and 0.795 for lysozyme. This suggests that the  
28 description of Ref. 18 fits better the lysozyme case than the BPTI case. It should also be  
29 noted that, since the PC-modes involved in these comparisons are already converged, the  
30 values obtained for the RMSIP<sub>2</sub> are insensitive to the total simulation time. To probe this  
31 we run extra trajectories of 5-ns, so that we had a set of 400 of such trajectories. Then  
32 we concatenated the 400 trajectories and calculated the PC-modes. These PC-modes were  
33 compared with those obtained by concatenating 40 trajectories of 50 ns so that, in the two  
34 cases, the total simulation time was 2.0  $\mu$ s. In this case, the value of RMSIP<sub>2</sub> was 0.603 for  
35 BPTI and 0.804 for lysozyme, which are nearly the same as reported above. It could also  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

happen that, in the extra time, some trajectories move between different wells. This would be rare for stable conformations of proteins if using the simulation times typically employed nowadays. However, the chances of observing such transitions increase as the number of concatenated trajectories gets large.<sup>33</sup> In that case, a more drastic effect on the converged PC-modes should be observed. As stated above, the method described in this article is not meant to address those cases. Nevertheless, it could be used to characterize the fluctuations observed within each well after the conformations sampled in the simulations have been clusterized so that they can be ascribed to each well. In any case, the most important conclusion to be drawn from this discussion is that PC-modes obtained by concatenating  $n$  trajectories of a simulation time  $T_s$  are representative of the configurational space attainable in time  $T_s$ . In general, they are not expected to be the same as PC-modes computed by concatenating longer trajectories.

## Conclusions

We have shown that concatenating  $n$  independent but equivalent MD simulations, and computing the PC-modes from the correlation matrix of the concatenated trajectory  $\mathbf{C}^{(n)}$ , significantly improves the reproducibility of the main PC-modes. The procedure has two important and convenient properties. First, small values of  $n$  provide a significant enhancement against the results obtained from single MD simulations. In particular, the possibility of getting badly-defined essential spaces is greatly reduced. Second, if desired or needed, the quality of the results can be systematically improved by increasing  $n$ . The main limitation of the procedure has also been stated. The PC-modes so obtained are representative of a specific simulation time: the time of the individual trajectories. They are not expected to be the same as PC-modes converged from longer simulations. We believe that the procedure proposed here will be particularly useful in quantitative applications of PCA, such as

1  
2  
3 calculations of entropy or free energy, as well as for approximate methods that rely in the  
4 selection of appropriate coordinates to reduce the dimensionality of the system.  
5  
6  
7

## 8 9 **Acknowledgement**

10  
11 The authors thank CONICET (Project ID 11220130100260CO) for the financial support  
12 of this research. Financial and computational support from the Universidad Nacional de  
13 Quilmes (Project ID 1402/15) is also gratefully acknowledged.  
14  
15  
16  
17  
18  
19

## 20 21 **Supporting information**

22  
23 A procedure to assess the consistency of the essential space determined by concatenating  
24 a given number of equivalent MD simulations is presented as Supporting information. The  
25 proposed method is completely general and easy to implement. It would allow to the users  
26 to recognize if the calculations already performed are enough or, on the contrary, some extra  
27 calculations are required. This information is available free of charge via the Internet at  
28 <http://pubs.acs.org> .  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## References

- (1) Amadei, A.; Linssen, A.; Berendsen, H. Essential Dynamics of Proteins. *Proteins: Struct., Funct., Bioinf.* **1994**, *17*, 412–425.
- (2) de Groot, B. L.; Amadei, A.; Scheek, R. M.; van Nuland, N. A. J.; Berendsen, H. J. C. An Extended Sampling of the Configurational Space of HPr from *E. coli*. *Proteins: Struct., Funct., Bioinf.* **1996**, *26*, 314–322.
- (3) Nymeyer, H.; García, A. E. Simulation of the Folding Equilibrium of  $\alpha$ -helical Peptides: a Comparison of the Generalized Born Approximation with Explicit Solvent. *Proc. Natl. Acad. Sci.* **2003**, *100*, 13934–13939.
- (4) Zhang, W.; Wu, C.; Duan, Y. Convergence of Replica Exchange Molecular Dynamics. *J. Chem. Phys* **2005**, *123*, 154105.
- (5) Karplus, M.; Kushick, J. N. Method for Estimating the Configurational Entropy of Macromolecules. *Macromolecules* **1981**, *14*, 325–332.
- (6) Perahia, D.; Levy, R.; Karplus, M. Motions of an  $\alpha$ -helical Polypeptide: Comparison of Molecular and Harmonic Dynamics. *Biopolymers* **1990**, *29*, 645–677.
- (7) Shenai, P. M.; Xu, Z.; Zhao, Y. Applications of Principal Component Analysis (PCA) in Materials Science. *Principal Component Analysis-Engineering Applications. InTech* **2012**, 25–40.
- (8) Yang, J. S.; Hu, X. J.; Li, X. X.; Li, H. Y.; Wang, Y. Application of Principal Component Analysis (PCA) for the Estimation of Source of Heavy Metal Contamination in Sediments of Xihe River, Shenyang City. *Advanced Materials Research*. 2012; pp 948–951.
- (9) Hayward, S.; Kitao, A.; Hirata, F.; Gō, N. Effect of Solvent on Collective Motions in Globular Protein. *J. Mol. Biol.* **1993**, *234*, 1207–1217.

- 1  
2  
3  
4 (10) Kitao, A.; Go, N. Investigating Protein Dynamics in Collective Coordinate Space. *Curr.*  
5 *Opin. Struct. Biol.* **1999**, *9*, 164–169.  
6  
7  
8 (11) Berendsen, H. J.; Hayward, S. Collective Protein Dynamics in Relation to Function.  
9 *Curr. Opin. Struct. Biol.* **2000**, *10*, 165–169.  
10  
11  
12 (12) Spoel, D. V. D.; de Groot, B. L.; Hayward, S.; Berendsen, H. J.; Vogel, H. J. Bending  
13 of the Calmodulin Central Helix: a Theoretical Study. *Protein Sci.* **1996**, *5*, 2044–2053.  
14  
15  
16 (13) van Aalten, D.; Jones, P.; De Sousa, M.; Findlay, J. Engineering Protein Mechanics:  
17 Inhibition of Concerted Motions of the Cellular Retinol Binding Protein by Site-directed  
18 Mutagenesis. *Protein Eng.* **1997**, *10*, 31–37.  
19  
20  
21 (14) Vesper, M. D.; de Groot, B. L.; Livesay, D. R. Collective Dynamics Underlying Al-  
22 losteric Transitions in Hemoglobin. *PLoS Comput. Biol.* **2013**, *9*, e1003232.  
23  
24  
25 (15) Hub, J. S.; de Groot, B. L. Detection of Functional Modes in Protein Dynamics. *PLoS*  
26 *Comput. Biol.* **2009**, *5*, e1000480.  
27  
28  
29 (16) Amadei, A.; Linssen, A. B.; de Groot, B. L.; Berendsen, H. J. Essential Degrees of  
30 Freedom of Proteins. **1995**, 85–93.  
31  
32  
33 (17) Balsera, M. A.; Wriggers, W.; Oono, Y.; Schulten, K. Principal Component Analysis  
34 and Long Time Protein Dynamics. *J. Phys. Chem.* **1996**, *100*, 2567–2572.  
35  
36  
37 (18) Caves, L. S.; Evanseck, J. D.; Karplus, M. Locally Accessible Conformations of Proteins:  
38 Multiple Molecular Dynamics Simulations of Crambin. *Protein Sci.* **1998**, *7*, 649–666.  
39  
40  
41 (19) Hess, B. Similarities between Principal Components of Protein Dynamics and Random  
42 Diffusion. *Phys. Rev. E* **2000**, *62*, 8438.  
43  
44  
45 (20) Amadei, A.; Ceruso, M. A.; Di Nola, A. On the Convergence of the Conformational Co-  
46 ordinates Basis Set Obtained by the Essential Dynamics Analysis of Proteins' Molecular  
47 Dynamics Simulations. *Proteins: Struct., Funct., Bioinf.* **1999**, *36*, 419–424.  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

- 1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60
- (21) Lambrughi, M.; Papaleo, E.; Testa, L.; Brocca, S.; De Gioia, L.; Grandori, R. Intramolecular Interactions Stabilizing Compact Conformations of the Intrinsically Disordered Kinase-inhibitor Domain of Sic1: a Molecular Dynamics Investigation. *Front Physiol.* **2012**, *3*, 435.
- (22) Martín-García, F.; Papaleo, E.; Gomez-Puertas, P.; Boomsma, W.; Lindorff-Larsen, K. Comparing Molecular Dynamics Force Fields in the Essential Subspace. *PLoS One* **2015**, *10*, e0121114.
- (23) Jónsdóttir, L. B.; Ellertsson, B. Ö.; Invernizzi, G.; Magnúsdóttir, M.; Thorbjarnardóttir, S. H.; Papaleo, E.; Kristjánsson, M. M. The Role of Salt Bridges on the Temperature Adaptation of Aqualysin I, a Thermostable Subtilisin-like Proteinase. *Biochim. Biophys. Acta* **2014**, *1844*, 2174–2181.
- (24) Pierdominici-Sottile, G.; Palma, J. New Insights into the Meaning and Usefulness of Principal Component Analysis of Concatenated Trajectories. *J. Comput. Chem* **2015**, *36*, 424–432.
- (25) Lou, H.; ; Cukier, R. I. Molecular Dynamics of Apo-Adenylate Kinase: A Distance Replica Exchange Method for the Free Energy of Conformational Fluctuations. *J. Phys. Chem. B* **2006**, *110*, 24121–24137.
- (26) Daidone, I.; Amadei, A. Essential Dynamics: Foundation and Applications. *Wiley Interdiscip. Rev. Comput. Mol. Sci.* **2012**, *2*, 762–770.
- (27) de Groot, B.; van Aalten, D.; Amadei, A.; Berendsen, H. The Consistency of Large Concerted Motions in Proteins in Molecular Dynamics Simulations. *Biophys. J* **1996**, *71*, 1707.
- (28) de Groot, B.; Hayward, S.; Van Aalten, D.; Amadei, A.; Berendsen, H. Domain Motions in Bacteriophage T4 Lysozyme: a Comparison between Molecular Dynamics and Crystallographic Data. *Proteins: Struct., Funct., Genet.* **1998**, *31*, 116–127.

- 1  
2  
3  
4 (29) Amadei, A.; Ceruso, M. A.; Di Nola, A. On the Convergence of the Conformational Co-  
5 ordinates Basis Set Obtained by the Essential Dynamics Analysis of Proteins' Molecular  
6 Dynamics Simulations. *Proteins: Struct., Funct., Bioinf.* **1999**, *36*, 419–424.  
7  
8  
9  
10 (30) Skjaerven, L.; Martínez, A.; Reuter, N. Principal Component and Normal Mode Analy-  
11 sis of Proteins; a Quantitative Comparison Using the GroEL Subunit. *Proteins: Struct.,*  
12 *Funct., Bioinf.* **2011**, *79*, 232–243.  
13  
14  
15  
16  
17 (31) Grossfield, A.; Feller, S. E.; Pitman, M. C. Convergence of Molecular Dynamics Simu-  
18 lations of Membrane Proteins. *Proteins: Struct., Funct., Bioinf.* **2007**, *67*, 31–40.  
19  
20  
21  
22 (32) van Aalten, D. M. F.; Amadei, A.; Linssen, A. B. M.; Eijssink, V. G. H.; Vriend, G.;  
23 Berendsen, H. J. C. The Essential Dynamics of Thermolysin: Confirmation of the  
24 Hinge-bending Motion and Comparison of Simulations in Vacuum and Water. *Proteins:*  
25 *Struct., Funct., Bioinf.* **1995**, *22*, 45–54.  
26  
27  
28  
29  
30  
31 (33) Pande, V. S.; Baker, I.; Chapman, J.; Elmer, S. P.; Khaliq, S.; Larson, S. M.;  
32 Rhee, Y. M.; Shirts, M. R.; Snow, C. D.; Sorin, E. J.; Zagrovic, B. Atomistic Protein  
33 Folding Simulations on the Submillisecond Time Scale Using Worldwide Distributed  
34 Computing. *Biopolymers* **2003**, *68*, 91–109.  
35  
36  
37  
38  
39  
40 (34) Ikebe, J.; Umezawa, K.; Kamiya, N.; Sugihara, T.; Yonezawa, Y.; Takano, Y.; Naka-  
41 mura, H.; Higo, J. Theory for Trivial Trajectory Parallelization of Multicanonical Molec-  
42 ular Dynamics and Application to a Polypeptide in Water. *J. Comput. Chem.* **2011**,  
43 *32*, 1286–1297.  
44  
45  
46  
47  
48  
49 (35) Allen, L. R.; Krivov, S. V.; Paci, E. Analysis of the Free-Energy Surface of Proteins  
50 from Reversible Folding Simulations. *PLoS Comput. Biol.* **2009**, *5*, e1000428.  
51  
52  
53  
54 (36) Hori, N.; Chikenji, G.; Berry, R. S.; Takada, S. Folding Energy Landscape and Network  
55 Dynamics of Small Globular Proteins. *Proc. Natl. Acad. Sci.* **2009**, *106*, 73–78.  
56  
57  
58  
59  
60



- 1  
2  
3  
4 (37) Sittel, F.; Jain, A.; Stock, G. Principal Component Analysis of Molecular Dynamics:  
5 On the Use of Cartesian vs. Internal Coordinates. *J. Chem. Phys.* **2014**, *141*.  
6  
7  
8  
9 (38) Faber, N.; Meinders, M.; Geladi, P.; Sjöström, M.; Buydens, L.; Kateman, G. Ran-  
10 dom Error Bias in Principal Component Analysis. Part I. Derivation of Theoretical  
11 Predictions. *Anal. Chim. Acta.* **1995**, *304*, 257–271.  
12  
13  
14  
15 (39) Nadler, B. Finite Sample Approximation Results for Principal Component Analysis: A  
16 Matrix Perturbation Approach. *Ann. Stat* **2008**, 2791–2817.  
17  
18  
19  
20 (40) Allez, R.; Bouchaud, J.-P. Eigenvector Dynamics: General Theory and some Applica-  
21 tions. *Phys. Rev. E* **2012**, *86*, 046202.  
22  
23  
24  
25 (41) Wlodawer, A.; Walter, J.; Huber, R.; Sjölin, L. Structure of Bovine Pancreatic Trypsin  
26 Inhibitor: Results of Joint Neutron and X-ray Refinement of Crystal form II. *J. Mol.*  
27 *Biol.* **1984**, *180*, 301–329.  
28  
29  
30  
31  
32 (42) Muraki, M.; Harata, K.; Sugita, N.; Sato, K. Origin of Carbohydrate Recognition  
33 Specificity of Human Lysozyme Revealed by Affinity Labeling,. *Biochemistry* **1996**,  
34 *35*, 13562–13567.  
35  
36  
37  
38  
39 (43) Hess, B. Convergence of Sampling in Protein Simulations. *Phys. Rev. E* **2002**, *65*,  
40 031910.  
41  
42  
43  
44 (44) Boechi, L.; de Oliveira, C. A. F.; Da Fonseca, I.; Kizjakina, K.; Sobrado, P.; Tan-  
45 ner, J. J.; McCammon, J. A. Substrate-Dependent Dynamics of UDP-galactopyranose  
46 Mutase: Implications for Drug Design. *Protein Sci.* **2013**, *22*, 1490–1501.  
47  
48  
49  
50  
51 (45) Lange, O. F.; Grubmüller, H. Full Correlation Analysis of Conformational Protein  
52 Dynamics. *Proteins: Struct., Funct., Bioinf.* **2008**, *70*, 1294–1312.  
53  
54  
55  
56 (46) Pierce, L. C.; Salomon-Ferrer, R.; de Oliveira, C. A. F.; McCammon, J. A.;  
57  
58  
59  
60

- 1  
2  
3 Walker, R. C. Routine Access to Millisecond Time Scale Events with Accelerated Molec-  
4 ular Dynamics. *J. Chem. Theory Comput.* **2012**, *8*, 2997–3002.  
5  
6  
7  
8  
9 (47) Lange, O. F.; Grubmüller, H. Generalized Correlation for Biomolecular Dynamics. *Pro-*  
10 *teins: Struct., Funct., Bioinf.* **2006**, *62*, 1053–1061.  
11  
12  
13 (48) Elofsson, A.; Nilsson, L. How Consistent are Molecular Dynamics Simulations? *J. Mol.*  
14 *Biol.* **1993**, *233*, 766 – 780.  
15  
16  
17  
18 (49) Schulze, B. G.; Evanseck, J. D. Cooperative Role of Arg45 and His64 in the Spec-  
19 *troscopic A3 state of Carbonmonoxy Myoglobin: Molecular Dynamics Simulations,*  
20 *Multivariate Analysis, and Quantum Mechanical Computations. J. Am. Chem. Soc.*  
21 **1999**, *121*, 6444–6454.  
22  
23  
24  
25  
26  
27 (50) Schiøtt, B.; Bruice, T. C. Reaction Mechanism of Soluble Epoxide Hydrolase: Insights  
28 from Molecular Dynamics Simulations. *J. Am. Chem. Soc.* **2002**, *124*, 14558–14570.  
29  
30  
31  
32 (51) Vitkup, D.; Ringe, D.; Karplus, M.; Petsko, G. A. Why Protein R-factors are so Large:  
33 A Self-consistent Analysis. *Proteins: Struct., Funct., Bioinf.* **2002**, *46*, 345–354.  
34  
35  
36  
37 (52) Worth, G. A.; Nardi, F.; Wade, R. C. Use of Multiple Molecular Dynamics Trajectories  
38 to Study Biomolecules in Solution: the YTGP Peptide. *J. Phys. Chem. B* **1998**, *102*,  
39 6260–6272.  
40  
41  
42  
43 (53) Gorfe, A. A.; Ferrara, P.; Caffisch, A.; Marti, D. N.; Bosshard, H. R.; Jelesarov, I.  
44 Calculation of Protein Ionization Equilibria with Conformational Sampling: pKa of a  
45 Model Leucine Zipper, GCN4 and Barnase. *Proteins: Struct., Funct., Bioinf.* **2002**,  
46 *46*, 41–60.  
47  
48  
49  
50  
51  
52  
53 (54) Perez, J. J.; Tomas, M. S.; Rubio-Martínez, J. Assessment of the Sampling Performance  
54 of Multiple-Copy Dynamics versus a Unique Trajectory. *J. Chem. Inf. Model.* **2016**,  
55 *56*, 1950–1962.  
56  
57  
58  
59  
60

- 1  
2  
3  
4 (55) Coveney, P. V.; Wan, S. On the Calculation of Equilibrium Thermodynamic Properties  
5 from Molecular Dynamics. *Phys. Chem. Chem. Phys.* **2016**, *18*, 30236–30240.  
6  
7  
8 (56) Shaw, D. E.; Maragakis, P.; Lindorff-Larsen, K.; Piana, S.; Dror, R. O.; East-  
9 wood, M. P.; Bank, J. A.; Jumper, J. M.; Salmon, J. K.; Shan, Y.; Wriggers, W.  
10 Atomic-Level Characterization of the Structural Dynamics of Proteins. *Science* **2010**,  
11 *330*, 341–346.  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## Figures

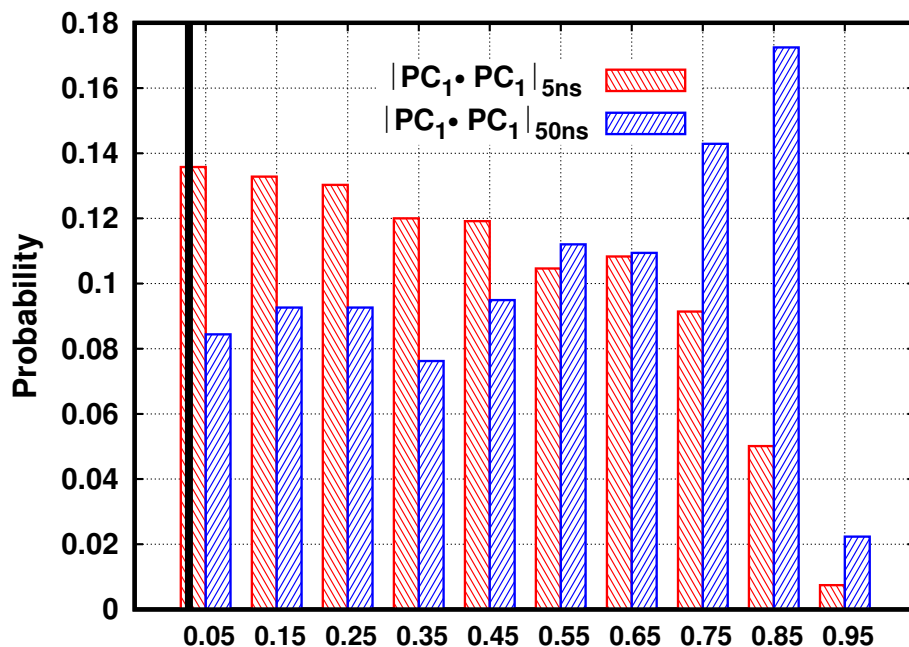


Figure 1: Normalized probability distributions for  $|\text{PC}_1^a \cdot \text{PC}_1^b|$ , computed from independent MD simulations of 5 ns (red) and 50 ns (blue) of BPTI. The vertical black line indicates the value of the inner product that contains 99% of cumulative probability, for normalized random vectors of the same dimension (see text).

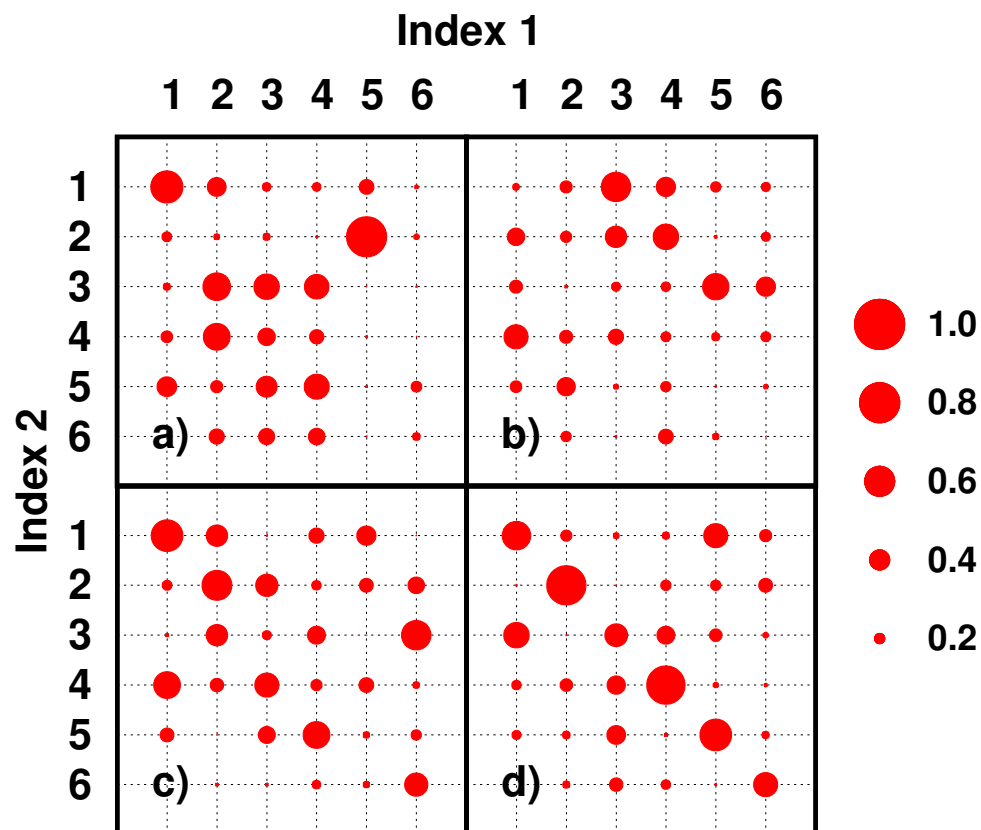


Figure 2: Typical examples of inner-product matrices for PC-modes computed from two independent MD simulations of BPTI. (a) and (b) correspond to 5-ns trajectories; (c) and (d) to 50-ns trajectories. The label of the axis refer to the index of the PC-modes while the radius of the circles measures the absolute value of the inner products.

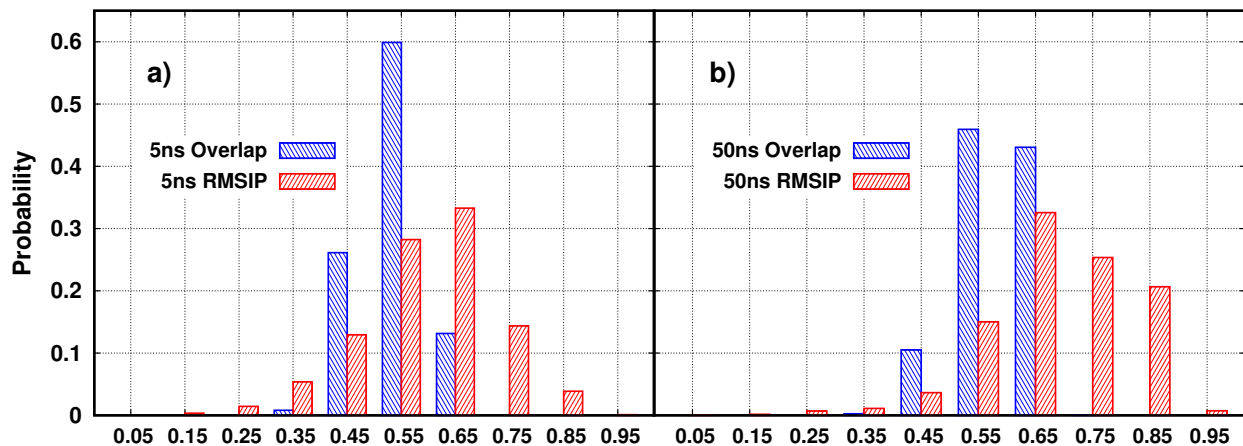


Figure 3: (a) Normalized probability distributions for  $\text{RMSIP}_2$  and overlap computed from the 16110 pairs of PC-mode sets formed from the 180 individual MD simulations of 5 ns of BPTI. (b) Normalized probability distributions for  $\text{RMSIP}_2$  and overlap computed from the 3160 pairs of PC-mode sets formed from the 80 individual MD simulations of 50 ns of BPTI.

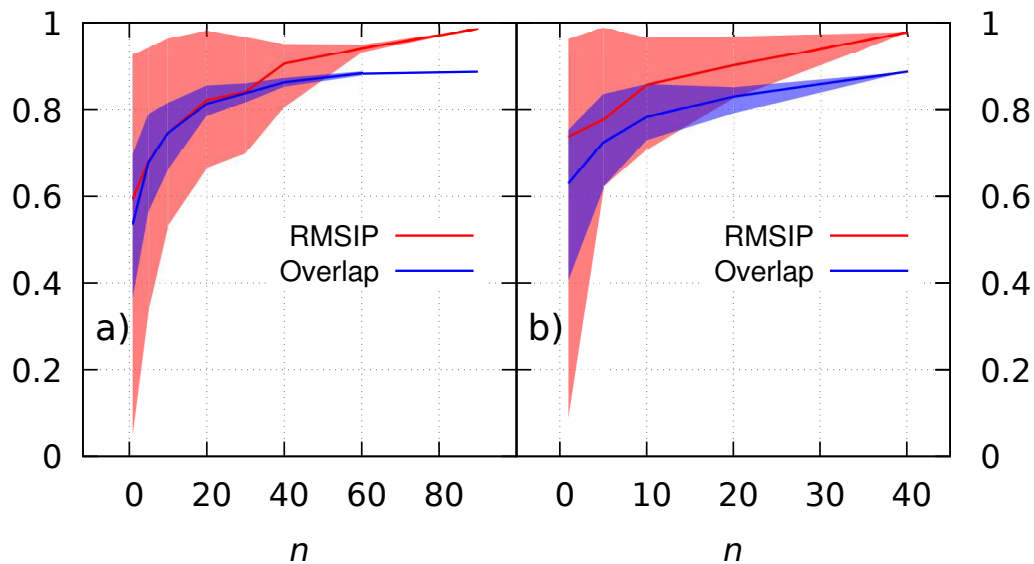


Figure 4: Evolution of  $\text{RMSIP}_2$  and overlap with the number of trajectories,  $n$ , included in the concatenated correlation matrix. Data correspond to simulations of BPTI. The solid lines indicate the average values. The shadows go from the minimum to the maximum value observed in the sample. (a) 5-ns trajectories; (b) 50-ns trajectories.

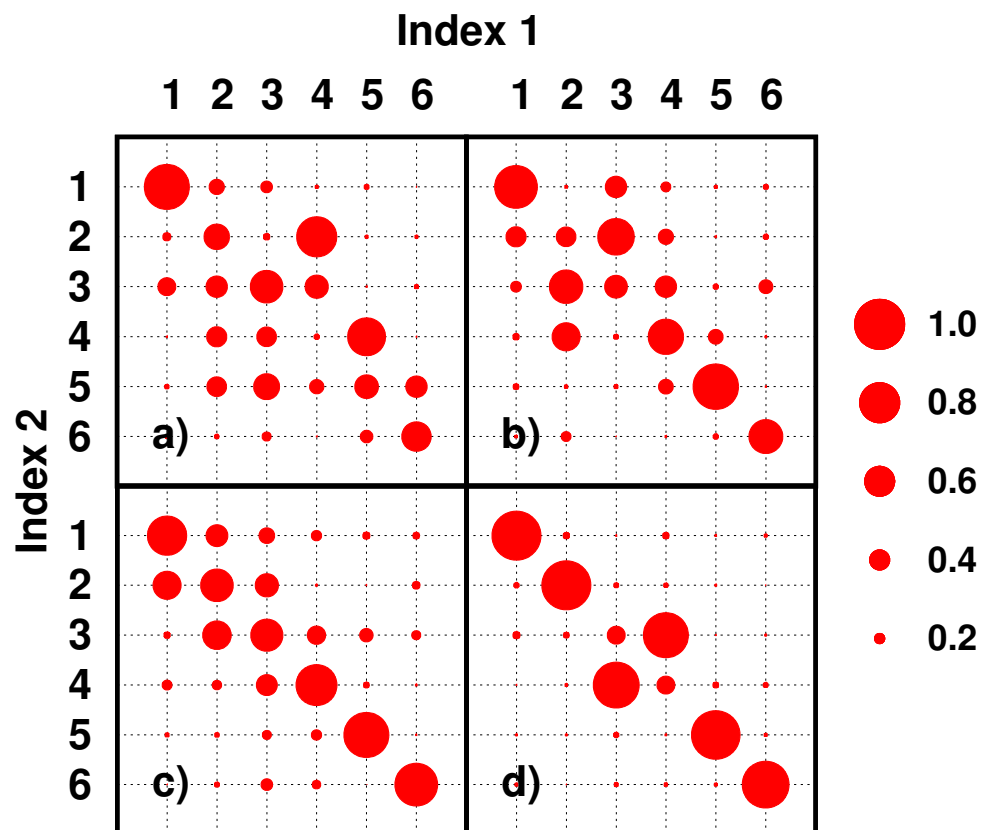


Figure 5: Inner product matrices for PC-modes computed from alternative batches of concatenated trajectories of BPTI. The label of the axis refer to the index of the PC-modes while the radius of the circles measures the absolute value of the inner products. a)  $n = 10$ ; b)  $n = 20$ ; c)  $n = 30$ ; d)  $n = 90$ .



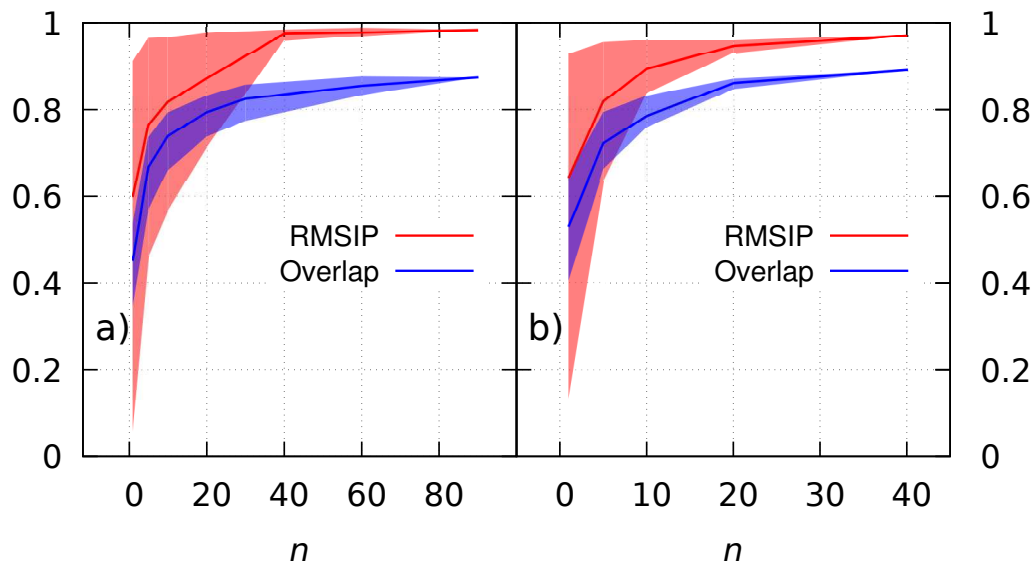


Figure 6: Evolution of  $\text{RMSIP}_2$  and overlap with the number of trajectories,  $n$ , included in the concatenated correlation matrix. Data correspond to simulations of lysozyme. The solid lines indicate the average values. The shadows go from the minimum to the maximum value observed in the sample. (a) 5-ns trajectories; (b) 50-ns trajectories.

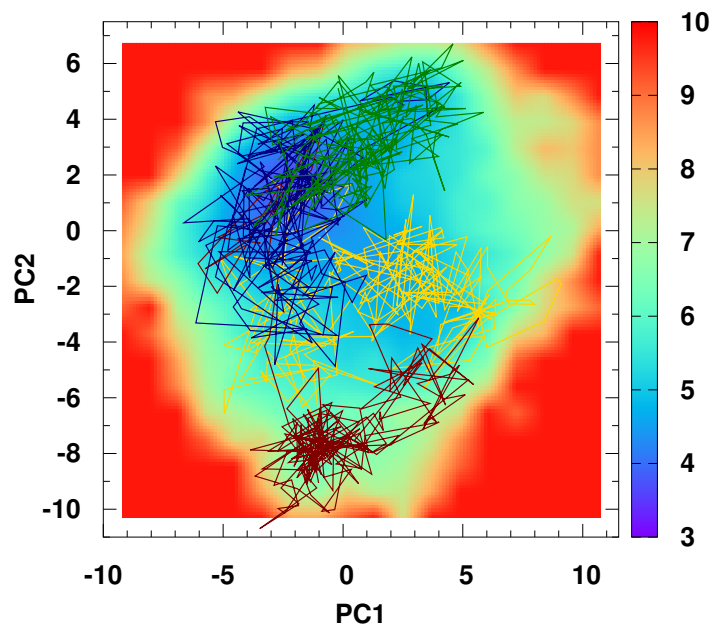


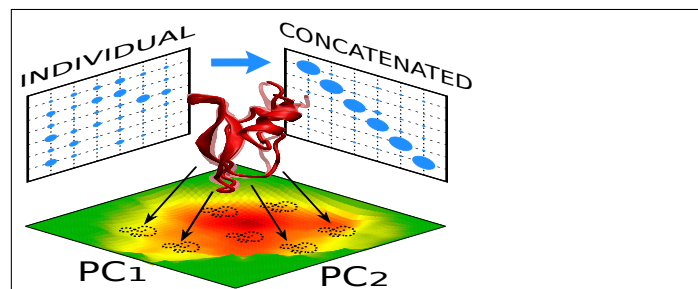
Figure 7: Typical individual trajectories projected onto the plane spanned by  $PC_1$  and  $PC_2$  of matrix  $\mathbf{C}^{(180)}$  of BPTI. The colored contour plot shows the free energy (in arbitrary units) computed from the snapshots collected along the 180 independent trajectories.

## Tables

Table 1: Number of batches ( $N_{batch}$ ) and number of pairs of batches ( $N_{pairs}$ ) that can be formed from the individual trajectories for each selected value of  $n$ .

180 traj (5 ns)			80 traj (50 ns)		
$n$	$N_{batch}$	$N_{pairs}$	$n$	$N_{batch}$	$N_{pairs}$
1	180	16110	1	80	3160
5	36	630	5	16	120
10	18	153	10	8	28
20	9	36	20	4	6
30	6	15	40	2	1
60	3	3			
90	2	1			

## Graphical TOC Entry

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60