Sampling Molecular Conformers in Solution with Quantum Mechanical Accuracy at a Nearly Molecular-Mechanics Cost

note finali coverpage

(Article begins on next page)

02 May 2024

# Sampling molecular conformers in solution with quantum mechanical accuracy at a nearly molecular mechanics cost

Marta Rosa[1], Marco Micciarelli[1], Alessandro Laio[1], Stefano Baroni[1,*]

**1 SISSA – Scuola Internazionale Superiore di Studi Avanzati
via Bonomea 265, 34136 Trieste – Italy**

**\* baroni@sissa.it**

## Abstract

We introduce a method to evaluate the relative populations of different conformers of molecular species in solution, aiming at quantum mechanical accuracy, while keeping the computational cost at a nearly molecular-mechanics level. This goal is achieved by combining long classical molecular-dynamics simulations to sample the free-energy landscape of the system, advanced clustering techniques to identify the most relevant conformers, and thermodynamic perturbation theory to correct the resulting populations, using quantum-mechanical energies from density-functional theory. A quantitative criterion for assessing the accuracy thus achieved is proposed. The resulting methodology is demonstrated in the specific case of cyanin (cyanidin-3-glucoside) in water solution.

# 1   Introduction

Organic molecules in water solution usually exist in several conformations, separated from each other by high free-energy barriers. Determining the relative population of these conformers is key for predicting molecular properties such as, *e.g.*, optical or NMR spectra. The procedure that is normally followed is to estimate these populations starting from the relative energies of the conformers, without treating explicitly the solvent molecules. The effect of the solvent can then be accounted for by using an implicit solvation scheme [Tomasi et al.(2005)Tomasi, Mennucci, and Cammi, Mennucci and Tomasi(1997)Mennucci, and Tomasi, Dupont et al.(2013)Dupont, Andreussi, and Marzari, Barone et al.(1997)Barone, Cossi, and Tomasi, Andreussi et al.(2012)Andreussi, Dabo, and Marzari, Orozco and Luque(2000)Orozco, and Luque] while entropic effects can be estimated in the harmonic approximation from the vibrational frequencies of the solute [Andricioaei and Karplus(2001)Andricioaei, and Karplus, Carlsson and Åqvist(2006)Carlsson, and Åqvist, Simonson et al.(2002)Simonson, Archontis, and Karplus, Suárez and Díaz(2015)Suárez, and Díaz]. This procedure is computationally expedient and provides an estimate of the populations that can at times be rather accurate. However, in many cases the procedure is affected by large systematic errors, due to the intrinsically molecular nature of the solvent. For example, a specific conformer can be stabilized by the presence of a solvent molecule bridging two moieties of the solute, a situation that would be missed by any implicit-solvent scheme, or the presence of floppy vibrational modes could make the use of the harmonic approximation questionable.

In this paper we propose an approach that allows estimating the relative populations of various conformers with the accuracy of ab initio (AI) molecular dynamics (MD) in explicit solvent at the cost of a few thousand quantum calculations. The configurational space of the solvated molecule is sampled by long molecular-mechanics (MM) MD runs, while density functional theory (DFT) calculations are performed only on a carefully selected set of configurations. In our approach a recently proposed clustering algorithm [Rodriguez and Laio(2014)Rodriguez, and Laio] is first applied to a long, supposedly ergodic, MMMD trajectory, to identify molecular conformers corresponding to suitably defined slow collective variables, and the relative populations are then estimated from the resulting residence times. In order to estimate quantum mechanical (QM) corrections to the population of these conformers, we exploit first-order thermodynamic perturbation theory, a procedure first pioneered by Warshel [Muller and Warshel(1995)Muller, and Warshel, Wesolowski and Warshel(1994)Wesolowski, and Warshel], using the QM energies computed at the DFT level on a set of uncorrelated configurations for each conformer. A key prerequisite for the success of this procedure is a high level of consistency between the MM and the QM free-energy landscapes. In particular, the stable conformers must be structurally similar at the MM and QM levels, while their populations can differ substantially. Importantly, we show that the magnitude of the second-order corrections to the conformational free energies, computed without taking into account the solvent, while not used directly to evaluate populations, provides a fair criterion to appraise *ex-post* the reliability of first-order corrections. Our approach is demonstrated by applying it to the cyanin (cyanidin-3-glucoside) molecule in water solution at room temperature.

# 2   Methodology

The configurational space of large molecules in solution is made of several free-energy basins, (*molecular conformers*) separated by high free-energy barriers. These systems can be easily described with MM simulations, and the relative populations of different

conformers estimated from the residence times of the molecule in each of them. Nevertheless, quantum accuracy in the inter-atomic forces is often needed, and sampling the configuration space with AIMD is hindered by the long time needed for the molecule to hop between any two conformers.

Here we describe an approach that allows one to estimate the relative populations of different conformers and to determine statistical averages of various observables at the QM level from MMMD trajectories. In the particular case where one is just interested in the relative populations of different conformers, the following derivation will lead to a direct estimate of the relevant free energy differences (see Eqs. 6 and 7).

In classical statistical mechanics, the physical properties of a system can be expressed as time averages of suitably defined configuration-space functions (*observables*) over molecular trajectories:

$$\langle A \rangle = \lim_{\tau \to \infty} \frac{1}{\tau} \int_0^\tau A(q_t) dt, \tag{1}$$

$$= \frac{1}{Z} \int A(q) e^{-\frac{V(q)}{k_B T}} dq, \tag{2}$$

where $q$ indicates the set of atomic coordinates, $V(q)$ is the potential energy function, $\tau$ the length of the trajectory, $T$ the system's temperature, $k_B$ the Boltzmann's constant, and $Z$ the partition function. Molecular conformers in configurational space are defined in such a way that thermalization within each of them occurs in typical molecular-vibration times, whereas the hopping between any two of them is kinetically hindered by free energy barriers: the corresponding transition rate follows therefore an Arrhenius law, typical of thermally activated processes. As a consequence, global thermalization requires simulation times that are exponentially long in the height of the free-energy barriers, and can thus be hardly achieved using accurate QM methods, whose scope is limited to processes spanning hundreds of pico-seconds at most. Methods based on MMMD extend by several orders of magnitude the time scales accessible to molecular simulations, bringing them up to the micro-second range and up, but their accuracy is limited by the quality of the empirical force fields employed therein.

In those cases where classical force fields are accurate enough to describe the general topography of the free-energy landscape, the computational convenience of MMMD can actually be combined with the accuracy of AIMD, by means of thermodynamic perturbation theory. In order to proceed, let us start from Eq. (2) and rewrite it as:

$$\langle A \rangle \approx \sum_{\mathcal{C}} \frac{Z^{\mathcal{C}}}{Z} \frac{1}{Z^{\mathcal{C}}} \int_{q \in \mathcal{C}} A(q) e^{-\frac{V(q)}{k_B T}} dq, \tag{3}$$

$$\approx \sum_{\mathcal{C}} p^{\mathcal{C}} \frac{1}{\tau^{\mathcal{C}}} \int_0^{\tau^{\mathcal{C}}} A(q_t) dt, \tag{4}$$

where the sums extend over all the conformers $\mathcal{C}$; in Eq. (3) the integrals are restricted to the portions of configurational space characterizing a given conformer and $Z^{\mathcal{C}}$ is the corresponding restricted partition function; in Eq. (4) $p^{\mathcal{C}} = \frac{Z^{\mathcal{C}}}{Z} \approx \frac{\tau^{\mathcal{C}}}{\tau}$ is the probability that the system is found in conformer $\mathcal{C}$ and the time averages are evaluated over times $\tau^{\mathcal{C}}$ longer than the local thermalization time, but significantly shorter than the typical residence time within that same conformer. These restricted averages are accessible to AIMD, whereas evaluating $p^{\mathcal{C}}$ would require long trajectories that can only be generated via MMMD. In order to overcome this difficulty, we express these probabilities in terms of the conformers' free energies, $F^{\mathcal{C}}$, defined as:

$$p^{\mathcal{C}} \propto e^{-\frac{F^{\mathcal{C}}}{k_B T}}. \tag{5}$$

Eq. (5) allows us to express the ratio between the populations of a same conformer computed at different levels of theory in terms of the exponential of the corresponding free-energy differences, [Zwanzig(1954)] such as resulting from two different MM force fields or a force field and a QM approach, such as DFT. In the latter case, with obvious notation one would have:

$$F_{QM}^{\mathcal{C}} = F_{MM}^{\mathcal{C}} + k_B T \ \log \left\langle \mathrm{e}^{\frac{1}{k_B T}(E_{QM}-E_{MM})} \right\rangle_{MM}^{\mathcal{C}}, \tag{6}$$

where $\langle \cdot \rangle_{MM}^{\mathcal{C}}$ indicates a time (or canonical) average performed within the $\mathcal{C}$ conformer at the MM level. Eq. (6) is in principle exact, if the $\mathcal{C}$ conformer is well defined both at the MM and QM levels, which obviously implies that the two levels of theory are close enough. Even in this case the evaluation of the statistical average of the exponential in Eq. (6) is severely hampered by thermal fluctuations in the exponent, which scale as the system size when solvent molecules are explicitly accounted for (more on reducing the impact of solvent energy fluctuations in the following). What can be done in practice is expanding the logarithm in Eq. (6) in a series of cumulants, in the spirit of perturbation theory, and retaining only the linear term, which converges relatively easily in water:

$$F_{QM}^{\mathcal{C}} \approx F_{MM}^{\mathcal{C}} + \langle \Delta E \rangle + \mathcal{O}\left(\Delta E^2\right), \tag{7}$$

where $\Delta E = E_{QM} - E_{MM}$ is the difference between the QM and MM energies.

The second order correction to Eq. (7) is given by $\frac{\kappa_2}{2k_B T}$, where $\kappa_2 = \langle \Delta E^2 \rangle - \langle \Delta E \rangle^2$ is the second cumulant of $\Delta E$ and all the averages are sampled in the $\mathcal{C}$ conformer on the MM distribution of molecular configurations.

Second and higher-order cumulants converge worse and worse as their order increases. Already at second order, the energy fluctuations due to the solvent are too large to compute the second cumulant for systems comprising a few hundred solvent molecules. An option could be to account for the effects of a few solvent molecules that interact most strongly with (*i.e.* that are closest to) the solute. Beside the intrinsic arbitrariness of this procedure, its most important drawback would be the difficulty to estimate its accuracy. We decided therefore to stick to first order in Eq. (7) to correct the relative populations of different conformers, while second cumulants, as computed neglecting solvent effects (*i.e. in vacuo*), are used to estimate the accuracy of this first-order approximation, using the procedure outlined below.

The contribution of the second cumulants to the relative populations vanishes if their value is the same for the different conformers, as the relative free energies would not be affected by them in this case. The relative magnitude of the second cumulants for different conformers can however be reliably estimated by computing them for *dehydrated* molecular configurations, *i.e.* by comparing QM and MM energies corresponding to the molecular configurations generated by an explicit-solvent MMMD simulation, upon stripping off solvent molecules.* The cumulants thus obtained obviously cannot be used to improve our estimate of QM corrections to the relative populations of different conformers, but they do provide a fair estimate of the accuracy of the first-order approximation to this correction. A further piece of information that comes from the calculation of second cumulants is their convergence rate, *i.e.* their variance within a same conformer: when the MM and QM free-energy landscapes are similar, the value of the second cumulant of a conformer will easily converge, at least when neglecting solvent effects, while this will be increasingly difficult as the MM and QM landscapes differ from each other. All in all, second cumulants determine the accuracy of the first-order corrections to the relative populations in three distinct ways: *i)* The statistical error of the first cumulant within a same conformer is given by the
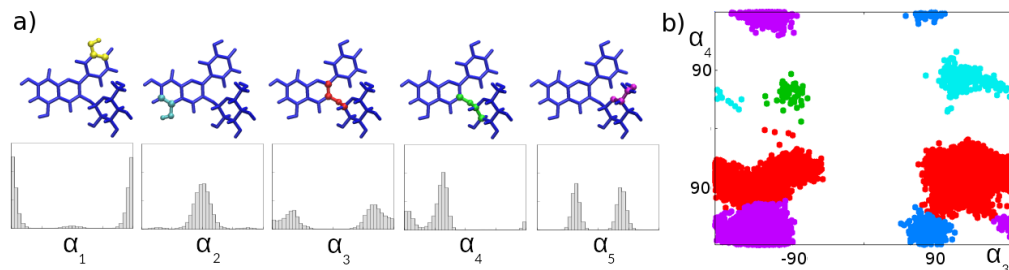
---

*Second cumulants can also be reliably estimated in a QM/MM scheme, because in this case energy fluctuations due to the solvent can be expediently removed relying on correlated sampling.

second cumulant divided by the number of independent molecular configurations generated in the MMMD run, $\Delta F_1^{\mathcal{C}} = \sqrt{\kappa_2^{\mathcal{C}}/N}$. *ii)* A second source of inaccuracy comes from the scatter of the value of the second cumulants across different conformers; a measure of this scatter is $\sigma_{\kappa_2}$, the standard deviation of the cluster distribution of second cumulants, and the corresponding contribution to the free-energy uncertainty is $\Delta F_2 = \frac{\sigma_{\kappa_2}}{2k_B T}$. *iii)* Finally, the second energy cumulant within each conformer is affected by its own statistical error, $\Delta \kappa_2^{\mathcal{C}}$, whose contribution to the overall uncertainty is $\Delta F_3^{\mathcal{C}} = \frac{\Delta \kappa_2^{\mathcal{C}}}{2k_B T}$. Combining these three contributions, the overall uncertainty on the free energy of the $\mathcal{C}$ conformer is estimated as:

$$\Delta F^{\mathcal{C}} = \sqrt{(\Delta F_1^{\mathcal{C}})^2 + (\Delta F_2)^2 + (\Delta F_3^{\mathcal{C}})^2}, \tag{8}$$

# 3  Results

As a case study, we tested our method on cyanin (cyanidin-3-glucoside, Figure 1). For this system we first ran a 2 $\mu$s MM trajectory using the Gromacs 4 MD package [Hess et al.(2008)Hess, Kutzner, Van Der Spoel, and Lindahl, Van Der Spoel et al.(2005)Van Der Spoel, Lindahl, Hess, Groenhof, Mark, and Berendsen]. The cyanin FF (referred to in the following as "FF$_{\text{amb}}$") was generated with the *antechamber* tool [Wang et al.(2006)Wang, Wang, Kollman, and Case, Wang et al.(2004)Wang, Wolf, Caldwell, Kollman, and Case] and RESP charges were calculated with the RESP ESP charge Derive (R.E.D.) program [Vanquelef et al.(2011)Vanquelef, Simon, Marquant, Garcia, Klimerak, Delepine, Cieplak, and Dupradeau]. The calculations were carried out in the NVT ensemble using a Nosé-Hoover thermostat in an orthorhombic cell with dimensions $19.22 \times 19.22 \times 16.53$ Å$^3$, filled with 177 water molecules. The slowest degrees of freedom for this molecule are dihedral rotations. Some of the relevant dihedrals are shown in Figure 1a together with their probability distribution evaluated along the MM trajectory. The distributions are peaked at a few maxima, signaling the presence of different conformers. The rotations of OH groups usually thermalize over AIMD time scales, with the exception of the $\alpha_1$ and $\alpha_2$ dihedrals depicted in Figure 1a.



**Figure 1.** Cyanin molecule in its neutral quinonoidal base state. (a) Upper row: the dihedrals highlighted in the different panels are those used to characterize the conformers. Lower row: probability distribution of the dihedrals depicted in the upper row, computed on a 2$\mu$s MM trajectory. (b): cluster representation of the conformational macrostates in the space of the $\alpha_3$ and $\alpha_4$ dihedrals; the different conformers found following Ref. [Rodriguez and Laio(2014)Rodriguez, and Laio] are highlighted in different colors.
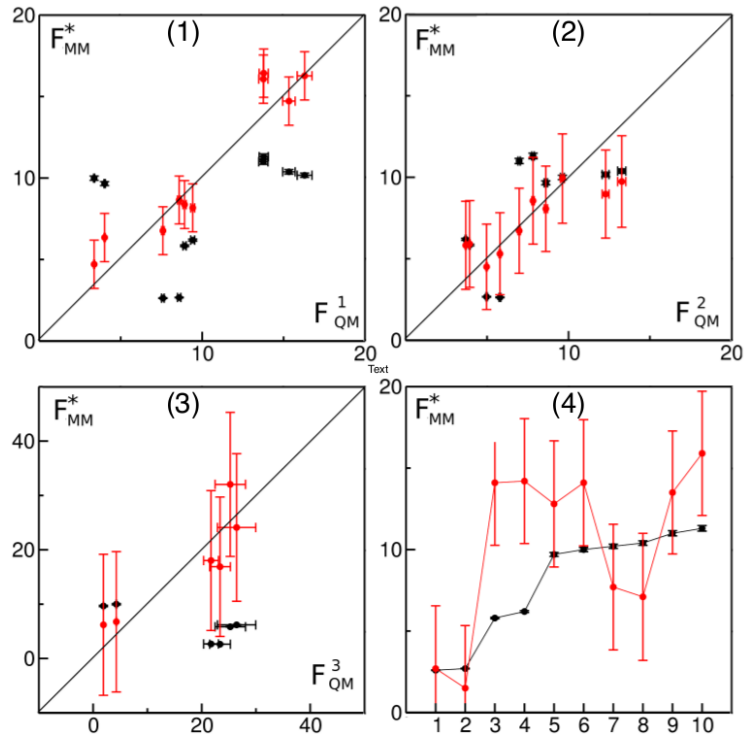
In Figure 1b we display the joint probability distributions of two different dihedrals. As typical in molecules of this complexity, different dihedrals are highly correlated, making the identification of the relevant conformers a non trivial task. In order to achieve this goal in the space defined by the five dihedrals of Figure 1a, we adopted a

newly developed clustering algorithm [Rodriguez and Laio(2014)Rodriguez, and Laio]. In this approach a cluster is defined as a peak in the joint probability density in the relevant space. The clusters identified by the algorithm in the subspace defined by the two dihedrals $\alpha_3$ and $\alpha_4$ are represented in different colors in Figure 1b. The density peaks are identified based on the observation that they have a larger density than neighboring configurations and that their minimum distance from points of higher density is anomalously large. In the case of water-solvated cyanin, we found 10 configurations satisfying this criterion in the space of the 5 dihedrals of Figure 1a. Once the density peaks have been thus identified, each of them is assumed to define a different conformer (cluster), and every other configuration is assigned to the cluster identified by the nearest density peak, following the procedure outlined in Ref. [ Rodriguez and Laio(2014)Rodriguez, and Laio]. The populations of these conformers according to the classical FF are defined as the number of configurations generated by the MMMD run and assigned to each of them by the clustering algorithm. Conformational populations cannot be evaluated in this way at the QM level of theory, because sufficiently long MD runs cannot be afforded in this case, and one has to resort therefore to a scheme based on free-energy corrections, following Eqs. (5) and (7).

A direct validation of this approach cannot be afforded either, because it would require a non-perturbative estimate of the QM populations, which is not feasible, as we have seen. For this reason, we decided to validate our perturbative approach using two different force fields, the first being the actual force field used in practice for our cluster analysis, $FF_{amb}$, the second representing a "classical proxy" of the QM description of the system and dubbed "$FF_{QM}$". We can then proceed to comparing the predictions of Eq. (7) for different proxies ($FF_{QM}^1, FF_{QM}^2, \cdots$), designed so as to represent increasingly important differences between the QM and MM levels of theory, against accurate estimates based on long MD runs performed with the same FFs. The goal of this procedure is to identify suitable indicators enabling us to estimate the accuracy of the perturbative approach when a non-perturbative evaluation of the relative populations (or free-energy differences) is not feasible.

We stress that in all cases molecular configurations are classified in terms of the clusters identified from the trajectory generated through the *bona fide* MM force field, $FF_{amb}$. "QM" conformational populations are estimated by assigning each configuration of the $FF_{QM}$ trajectory to the closest cluster as defined by clustering the $FF_{amb}$ trajectory. In this way $FF_{QM}$ populations can be estimated as the number of geometries belonging to each conformer, and the corresponding free energies evaluated from Eq. (5). The typical auto-correlation time of the energy within each conformer is of the order of 3-4 ps. In order to sample QM-MM energy differences, 2000 configurations per conformer were thus selected along the MM trajectory with a time lag of 10 ps. Second energy cumulants were also evaluated over this same set of configurations, once water molecules had been stripped off from them, thus permitting to considerably reduce statistical fluctuations, without compromising our aim to utilize second cumulants to estimate the accuracy of first-order free-energy differences.

In Figure 2 (a-c) we show the correlation plots between $FF_{QM}$ and $FF_{amb}$ free energies for three different QM proxy FFs, dubbed as $FF_{QM}^1$, $FF_{QM}^2$, and $FF_{QM}^3$, and designed so as to differ increasingly from the *bona fide* MM FF, $FF_{amb}$. In each one of these figures is reported the $FF_{QM}$ free energy (x axis) correlated with the $FF_{amb}$ free energy (black) and with the proxy QM free energy (red) as calculated from first-order thermodynamical perturbation theory from $FF_{amb}$. The errors on the linear corrections are estimated from the second energy cumulants, as explained at the end of the previous section. One clearly sees that the magnitude of the estimated errors is a faithful indicator of the quality of the first-order approximation: free energies calculated with Eq. (7) from $FF_{amb}$ results can always be trusted within the estimated error. By

**Figure 2.** Panels 1-3. Horizontal axis: $FF_{QM}$ is the "quantum" free energy of different molecular conformers computed with different "classical proxy" force fields ($FF_{QM}^1$, $FF_{QM}^2$, and $FF_{QM}^3$ in panels 1-3, respectively: see text). Vertical axis: $FF_{QM}^*$ is the free energy computed for each conformer using Eq. (7). Black symbols indicate purely "classical" results obtained with the $FF_{amb}$ force field and neglecting first order corrections in Eq. (7). Red symbols indicate results including first order corrections. Error bars are estimated from Eq. (8). Note the change of scale between panels 1-2 and 3. Panel 4: comparison of the classical and QM free energies (the latter computed at the DFT level using the perturbative scheme introduced in the present paper) of different molecular conformers of the cyanin quinoidal base. Units are kJ/mol throughout.

studying *a posteriori* the free energy landscapes generated by the $FF_{QM}^1$, $FF_{QM}^2$, and $FF_{QM}^3$ force fields, we notice how the corresponding configurational spaces become increasingly different from that of $FF_{amb}$. In the extreme case of $FF_{QM}^3$ the two systems are so different that the 10 most populated conformers of FF do not correspond anymore to the most populated conformers of $FF_{QM}^3$ (this is the reason why in Figure 2-c only six conformers are shown).

Finally, we applied our method to evaluate the QM free energies obtained at the DFT level. DFT calculations were performed with the QUANTUM ESPRESSO package version 5.0, [Giannozzi et al.(2009)Giannozzi, Baroni, Bonini, Calandra, Car, Cavazzoni, Ceresoli, Chiarotti, Cococcioni, Dabo, Dal Corso, de Gironcoli, Fabris, Fratesi, Gebauer, Gerstmann, Gougoussis, Kokalj, Lazzeri, Martin-Samos, Marzari, Mauri, Mazzarello, Paolini, Pasquarello, Paulatto, Sbraccia, Scandolo, Sclauzero, Seitsonen, Smogunov, Umari, and Wentzcovitch] with the same cell and number of solvent molecules as used in the MM simulations illustrated above, following the same procedure as with the proxy QM FFs. DFT estimates of the various free energies are shown in Figure 2-DFT.

The final value of the estimated error is very similar for the different clusters and is on average 3.8 kJ/mol. The magnitude of the error can be further reduced by increasing the size of the configuration sample and/or by designing a classical FF that more closely mimics the DFT inter-atomic interactions, so as to improve the quality of predictions based on first-order perturbation theory. We stress that the estimate of the error is also a function of the number of conformers we consider. In the present case, for instance, if we focus on the free energies of the conformers differing only for the orientation of the sugar (1, 2, 3, 4, 8, and 9) the error estimates lowers to $\approx 2.5$kJ/mol, showing a very good performance of the FF. Focusing on the pairs of conformers 1-2, 4-5, and 6-7, which differ for the orientation of $\alpha_1$ or $\alpha_2$ dihedrals (Figure 1 a), instead, the error estimate raises to 3.5 kJ/mol.

In order to get further insight into the free-energy differences predicted by our methodology, AIMD simulations were started from geometries belonging to conformers 1, 2, and 3. The latter trajectory was observed to move spontaneously towards a more stable conformer (1 or 2) after few ps of dynamics. This result shows that conformer 3 is unlikely to be stable at the QM level of theory, consistently with large free-energy differences between it and conformers 1-2 ($> 10$kJ/mol) and with a very low barrier in-between, which, if existing, is easily crossed at room temperature.

## 4   Conclusions

In this paper we have introduced a new method to sample and characterize the conformational space of complex molecular species in solution, using first-order thermodynamical perturbation theory to estimate quantum-mechanical corrections to classical molecular-dynamics results, and second-order perturbation theory to estimate the ensuing accuracy.

First-order perturbation theory has been widely used in the past to evaluate quantum-mechanical corrections to free-energy differences, as estimated from classical molecular dynamics, particularly in the study of chemical reactions and solvation free energies [Wesolowski and Warshel(1994)Wesolowski, and Warshel, Brandsdal et al.(2003)Brandsdal, Osterberg, Almlof, Feierberg, Luzhkov, and Aqvist, Hu and Yang(2008)Hu, and Yang, Rosta et al.(2006)Rosta, Klähn, and Warshel, Woods et al.(2008)Woods, Manby, and Mulholland]. The scope of this methodology when applied to complex molecular systems is limited by the ability of the low level of theory (MM in our case) to describe the zero-th order conformational landscape with sufficient accuracy. Whilst this description can be systematically improved, at least in principle, by improving the MM force field, no reliable criteria have been available so far to evaluate the quality of the first-order correction. One of the main steps forward made in our work is the identification of such a quantitative criterion, based on a careful analysis of an approximate evaluation of the second-order correction. Our analysis also indicates that first-order corrections need not be small in order to be accurate: the accuracy only depends on the ability of the low level of theory to correctly describe the topography of the conformers, which in turn can be assessed by the accuracy criterion mentioned above.

A second important element of this work is a novel procedure for identifying the conformers on which perturbation theory is applied. A conformer is defined as a peak in the probability density in the possibly high-dimensional space spanned by the internal coordinates of the solute. These peaks are identified by a novel clustering algorithm [Rodriguez and Laio(2014)Rodriguez, and Laio]. Since the analysis is performed at finite temperature and in the presence of a solvent, the conformers do not necessarily coincide with the minima of the potential energy surface. A correct definition of the conformers is a crucial ingredient in the procedure since, as we already

mentioned, an error in the description of the topography of the system dramatically impacts the accuracy.

Of course, the problem still remains that the perturbative evaluation of free-energy differences is hindered by statistical fluctuations, which crucially depend on the system size. In the present case size-extensive fluctuations due to an explicit account of the solvent can be significantly reduced by relying on a QM/MM approach whereby the extensive contributions of the solvent cancel exactly when computing the difference between QM/MM and pure-MM energies.

In short, we believe that our approach offers a route to probe the free energy landscape of highly mobile molecular species with a QM level of description *by keeping the statistical accuracy under control.* The availability of a reliable estimate of the error naturally opens the way to systematically improve the free-energy estimates by fine-tuning the low level of theory.

# References

Tomasi et al.(2005)Tomasi, Mennucci, and Cammi. Tomasi, J.; Mennucci, B.; Cammi, R. *Chemical Reviews* **2005**, *105*, 2999–3094.

Mennucci and Tomasi(1997)Mennucci, and Tomasi. Mennucci, B.; Tomasi, J. *The Journal of Chemical Physics* **1997**, *106*, 5151–5158.

Dupont et al.(2013)Dupont, Andreussi, and Marzari. Dupont, C.; Andreussi, O.; Marzari, N. *The Journal of Chemical Physics* **2013**, *139*, 214110.

Barone et al.(1997)Barone, Cossi, and Tomasi. Barone, V.; Cossi, M.; Tomasi, J. *The Journal of Chemical Physics* **1997**, *107*, 3210–3221.

Andreussi et al.(2012)Andreussi, Dabo, and Marzari. Andreussi, O.; Dabo, I.; Marzari, N. *The Journal of Chemical Physics* **2012**, *136*, 064102.

Orozco and Luque(2000)Orozco, and Luque. Orozco, M.; Luque, F. J. *Chemical Reviews* **2000**, *100*, 4187–4226.

Andricioaei and Karplus(2001)Andricioaei, and Karplus. Andricioaei, I.; Karplus, M. *The Journal of Chemical Physics* **2001**, *115*, 6289–6292.

Carlsson and Åqvist(2006)Carlsson, and Åqvist. Carlsson, J.; Åqvist, J. *Physical Chemistry Chemical Physics* **2006**, *8*, 5385–5395.

Simonson et al.(2002)Simonson, Archontis, and Karplus. Simonson, T.; Archontis, G.; Karplus, M. *Accounts of Chemical Research* **2002**, *35*, 430–437.

Suárez and Díaz(2015)Suárez, and Díaz. Suárez, D.; Díaz, N. *Wiley Interdisciplinary Reviews: Computational Molecular Science* **2015**, *5*, 1–26.

Rodriguez and Laio(2014)Rodriguez, and Laio. Rodriguez, A.; Laio, A. *Science* **2014**, *344*, 1492–1496.

Muller and Warshel(1995)Muller, and Warshel. Muller, R. P.; Warshel, A. *The Journal of Physical Chemistry* **1995**, *99*, 17516–17524.

Wesolowski and Warshel(1994)Wesolowski, and Warshel. Wesolowski, T.; Warshel, A. *The Journal of Physical Chemistry* **1994**, *98*, 5183–5187.

Zwanzig(1954). Zwanzig, R. W. *The Journal of Chemical Physics* **1954**, *22*, 1420–1426.

Hess et al.(2008)Hess, Kutzner, Van Der Spoel, and Lindahl. Hess, B.; Kutzner, C.; Van Der Spoel, D.; Lindahl, E. *Journal of Chemical Theory and Computation* **2008**, *4*, 435–447.

Van Der Spoel et al.(2005)Van Der Spoel, Lindahl, Hess, Groenhof, Mark, and Berendsen. Van Der Spoel, D.; Lindahl, E.; Hess, B.; Groenhof, G.; Mark, A. E.; Berendsen, H. J. *Journal of Computational Chemistry* **2005**, *26*, 1701–1718.

Wang et al.(2006)Wang, Wang, Kollman, and Case. Wang, J.; Wang, W.; Kollman, P. A.; Case, D. A. *Journal of Molecular Graphics and Modelling* **2006**, *25*, 247–260.

Wang et al.(2004)Wang, Wolf, Caldwell, Kollman, and Case. Wang, J.; Wolf, R. M.; Caldwell, J. W.; Kollman, P. A.; Case, D. A. *Journal of Computational Chemistry* **2004**, *25*, 1157–1174.

Vanquelef et al.(2011)Vanquelef, Simon, Marquant, Garcia, Klimerak, Delepine, Cieplak, and Vanquelef, E.; Simon, S.; Marquant, G.; Garcia, E.; Klimerak, G.; Delepine, J. C.; Cieplak, P.; Dupradeau, F.-Y. *Nucleic Acids Research* **2011**, *39*, W511–W517.

Giannozzi et al.(2009)Giannozzi, Baroni, Bonini, Calandra, Car, Cavazzoni, Ceresoli, Chiaro Giannozzi, P. et al. *Journal of Physics: Condensed Matter* **2009**, *21*, 395502 (19pp).

Brandsdal et al.(2003)Brandsdal, Osterberg, Almlof, Feierberg, Luzhkov, and Aqvist. Brandsdal, B. O.; Osterberg, F.; Almlof, M.; Feierberg, I.; Luzhkov, V. B.; Aqvist, J. *Advances in Protein Chemistry* **2003**, *66*, 123–158.

Hu and Yang(2008)Hu, and Yang. Hu, H.; Yang, W. *Annual Review of Physical Chemistry* **2008**, *59*, 573.

Rosta et al.(2006)Rosta, Klähn, and Warshel. Rosta, E.; Klähn, M.; Warshel, A. *The Journal of Physical Chemistry B* **2006**, *110*, 2934–2941.

Woods et al.(2008)Woods, Manby, and Mulholland. Woods, C. J.; Manby, F. R.; Mulholland, A. J. *The Journal of Chemical Physics* **2008**, *128*, 014109.