# Quantum harmonic free energies for biomolecules and nanomaterials

Alec F. White[1,2], Chenghan Li[1], Xing Zhang[1] and Garnet Kin-Lic Chan[1]

[1]Division of Chemistry and Chemical Engineering, California Institute of Technology, 1200 East California Boulevard, Pasadena, 91125, CA, United States of America.
[2]Present address: Quantum Simulation Technologies, Inc., Cambridge, MA 02139, United States of America.

**Abstract**

Obtaining the free energy of large molecules from quantum mechanical energy functions is a longstanding challenge. We describe a method that allows us to estimate, at the quantum mechanical level, the harmonic contributions to the thermodynamics of molecular systems of large size, with modest cost. Using this approach, we compute the vibrational thermodynamics of a series of diamond nanocrystals, and show that the error per atom decreases with system size in the limit of large systems. We further show that we can obtain the vibrational contributions to the binding free energies of prototypical protein-ligand complexes where the exact computation is too expensive to be practical. Our work raises the possibility of routine quantum mechanical estimates of thermodynamic quantities in complex systems.

The contributions to the free energy from atomic motion are critically important to the thermodynamics and kinetics of biological, chemical, and materials systems. Changes in such contributions govern processes ranging from the affinity of drug binding to structural phase transitions in crystals. When the internal energy is computed at the quantum mechanical level, a harmonic approximation is often the only feasible option to describe atomic motion. However, for large systems such as nanostructures and biomolecules, computing free energy contributions is expensive even within the harmonic approximation. For a system of $N$ atoms, the Hessian matrix which describes the vibrations requires $O(3N)$ gradient calculations, or $O(3N)$

1

times the cost of computing the internal energy. This is clearly prohibitive when $N$ is large, making free-energy computation in large systems with quantum mechanical methods a major contemporary challenge [1].

There are many possible strategies to speed up harmonic vibrational analysis, including methods based on a partial Hessian[2, 3], iterative diagonalization [4], and Hessian-free methods that use molecular dynamics to approximate the the harmonic problem[5, 6]. Here we describe a different strategy where we estimate vibrational thermodynamic quantities directly without ever computing the full Hessian or taking advantage of any local structure.

The starting point is to express each harmonic thermodynamic quantity as a matrix function trace. Then, our technique contains three elements. First, we sample the matrix trace operation using random vectors and stochastic Lanczos quadrature [7]. Second, we compute the Hessian-vector product at the same cost as the gradient from the difference of gradients at displaced geometries, bypassing the Hessian construction entirely. Third, we ameliorate the stochastic error, especially for free energy differences, through a form of correlated sampling. Related stochastic methods have been used for anharmonic corrections to the harmonic free energy[8, 9], as well as in stochastic electronic structure[10], but to our knowledge this is the first time these ideas have been brought to bear on the harmonic thermodynamic quantities themselves. As we demonstrate, this allows us to compute at the quantum mechanical level and with modest cost, vibrational free energy contributions for nanocrystals with more than 600 atoms, and free energy differences in protein-ligand complexes with more than 3000 atoms.

## Theory

We first express the harmonic thermochemical quantities as traces of matrix functions. In particular, we are interested in the zero-point energy (ZPE),

$$\text{ZPE} = \sum_I \frac{\omega_I}{2} = \sum_I \frac{\sqrt{\omega_I^2}}{2} = \text{Tr}\left[\frac{\sqrt{\mathbf{D}}}{2}\right] \tag{1}$$

the thermal contribution to the enthalpy,

$$H_{\text{vib}} - \text{ZPE} = \sum_I \omega_I \left(\frac{e^{-\beta\omega_I}}{1 - e^{-\beta\omega_I}}\right) = \text{Tr}\left[\frac{\sqrt{\mathbf{D}}\exp(-\beta\sqrt{\mathbf{D}})}{1 - \exp(-\beta\sqrt{\mathbf{D}})}\right] \tag{2}$$

and the entropy,

$$\begin{aligned} S_{\text{vib}}/k_B &= \sum_I \left[\beta\omega_I \frac{e^{-\beta\omega_I}}{1 - e^{-\beta\omega_I}} - \ln\left(1 - e^{-\beta\omega_I}\right)\right] \\ &= \text{Tr}\left[\beta\frac{\sqrt{\mathbf{D}}\exp(-\beta\sqrt{\mathbf{D}})}{1 - \exp(-\beta\sqrt{\mathbf{D}})} - \ln\left(1 - \exp(-\beta\sqrt{\mathbf{D}})\right)\right] \end{aligned} \tag{3}$$

Here $\beta$ is the inverse temperature, $\{\omega_I\}$ is the set of normal mode frequencies, and $\mathbf{D}$ is the mass-weighted Hessian matrix. We refer to ZPE as a non-thermal quantity as it has no temperature dependence.

The above expressions have the form $\mathrm{Tr}\, f(\mathbf{D})$. We now employ a stochastic estimator of the trace. The simplest version writes $\mathrm{Tr}\, f(\mathbf{D}) \approx \frac{M}{n} \sum_{l=1}^{n} \mathbf{v}_l^{\mathrm{T}} f(\mathbf{D})\mathbf{v}_l$ where $\mathbf{v}_l$ are a set of $n$ random vectors with zero mean and unit covariance, and $M = 3N - 6$ is the dimension of $\mathbf{D}$. This direct stochastic evaluation requires a polynomial approximation of $f(\mathbf{D})$, which is typically carried out using a Chebyshev expansion [11]. A closely related idea, which we use in this work, is stochastic Lanczos quadrature[7]. In this technique, the polynomial approximation is generated for the scalar $\mathbf{v}_l^T f(\mathbf{D})\mathbf{v}_l$ rather than globally for the function $f(\mathbf{D})$. We have found the stochastic Lanczos method to be slightly superior to the Chebyshev polynomial approach for the quantities in this work. Within the polynomial expansion, the main operation is the matrix-vector product $\mathbf{Dx}$, where $\mathbf{x}$ is in the stochastic Lanczos space. This can be computed from the difference of gradients, displaced by $\delta\mathbf{x}$ in mass-weighted coordinates, for small $\delta$ [12]. Thus no Hessian is needed at all in this approach (further details are provided in the Methods section). We note the sampling itself does not depend on $f$, and thus intermediate information, such as Lanczos quadrature weights and positions can be cached to estimate $\mathrm{Tr}\, f(\mathbf{D})$ for any other $f$, and as such thermostatistical quantities at different temperatures can be computed without the need for repeating the sampling procedure.

Within the above scheme, there are two sources of error. The first is from the order of the Lanczos quadrature, $m$. This vanishes when $m$ is greater than or equal to the matrix dimension. The second is the sampling error, which decreases like $1/\sqrt{n}$ for $n$ random vectors. For a quadrature order of $m$ and $n$ samples, the cost of the method is equal to $O(mn)$ gradient calculations. To reduce the statistical error, we use a form of correlated sampling. For the absolute thermodynamic quantities computed for the diamond nanocrystals we employ a high-level quantum mechanical approach as well as a cheaper low-level method (for example, a force-field, or semi-empirical quantum-mechanical approach) where the exact computation of the harmonic thermodynamic quantity $X_{\mathrm{low}}$ is possible. Then, the free energy contribution for the high-level method is obtained as

$$X_{\mathrm{high}} = X_{\mathrm{low}} + \Delta \tag{4}$$

where $\Delta$ is computed by applying the stochastic Lanczos quadrature to the difference of high-level and low-level methods. In the case of protein(P)-ligand(L) binding, we are interested in the difference between the holo (ligand-bound) state and the apo (ligand-free) state, i.e.

$$X_{\mathrm{bind}} = X_{\mathrm{P+L}} - X_{\mathrm{P}} - X_{\mathrm{L}} \tag{5}$$

where $X_{\mathrm{bind}}$ represents the ligand binding free energy, enthalpy, entropy, etc. In this case, we perform correlated sampling by using the same random vectors in the stochastic Lanczos treatment of the P+L, P, L systems (zeroing out elements for P and L respectively). No additional low-level method is involved in ligand-binding calculations.

# Results

As a first application, we take diamond nanocrystals (Figure 1a) as prototypical nano-materials, and compute the free energies as a function of size. We employ the high- and low-level correlated sampling approach described above, with Kohn-Sham density functional theory (DFT) with the PBE functional [13] as the high-level method and the semi-empirical extended tight-binding (xTB) method [14] as the low-level method. For the smallest system ($C_{54}H_{54}$), we can compute the Hessian explicitly to provide an exact reference. Figure 1b shows $\Delta$ quantities for the zero-point energy, the thermal enthalpy, and the entropy respectively, as a function of quadrature order $m$. The error bar indicates the statistical error for 50 samples (see SI for details of error analysis). A stochastic quadrature level of $m = 8$ does not provide sufficient accuracy, so we choose $m = 16$ for further calculations. Additional calculations on three transition-metal complexes using $m = 16$ are presented in the SI, although the choice of $m$ in general is system and accuracy specific. Figure 1c shows the value and stochastic error per atom (estimated as one standard error) for the ZPE, thermal enthalpy, and entropy respectively. We note the error decreases with the size of the system, faster than the decay of the quantities themselves, which is evidence of "self-averaging" due to the large system size. A more detailed discussion of the "self-averaging" behavior is provided in the SI.  Thus if one is interested in per-atom quantities, as is often the case for thermodynamics, for example to locate phase transitions, our stochastic approach becomes increasingly more efficient in a large system. The difference between our largest simulation and the extrapolated thermo-dynamic limits for the per-atom ZPE, enthalpy, and entropy is only 0.2 kcal/mol, 0.004 kcal/mol, and 0.006 kcal/mol respectively; statistical errors with 50 samples are about 0.001 kcal/mol or less. In fact, in the largest diamond system, with a *single* sample, one can estimate the per-atom quantities with a statistical error of less than 0.01 kcal/mol (this error was estimated from 50 samples; see SI for details) at a 120-fold speedup relative to the exact Hessian calculation.

If one is instead interested in the absolute values, the method can still be cheaper than the full computation of the Hessian. Figure 1d shows the stochastic error for the largest carbon system ($C_{432}H_{216}$) as a function of computational cost. At less than 20% of the exact calculation's computational effort, the error estimate for ZPE is well within 1 kcal/mol, red corresponding to 0.03% relative error with respect to the total ZPE or 7.7% relative error with respect to the $\Delta$ZPE between DFT and xTB. Less precise estimates can be obtained even more cheaply; a 20 times speedup is possible if 2 kcal/mol of error in the absolute quantities is tolerable.

Vibrational contributions to the free energy are also central to the study of large molecule and biomolecular interactions. For protein-ligand interactions in particular, where the aggregate binding is often only 5-15 kcal/mol, the vibrational contribution to binding free energies can be significant. Although the harmonic approximation is not necessarily a faithful approximation in these systems, the harmonic contributions nonetheless provide a useful first estimate of the thermal and entropic contribu-tions [15, 16]. The task is challenging for the stochastic Lanczos approach as binding is the difference of large absolute quantities, requiring tight convergence of the sta-tistical error. To reduce statistical error, we use the fixed random vector correlated
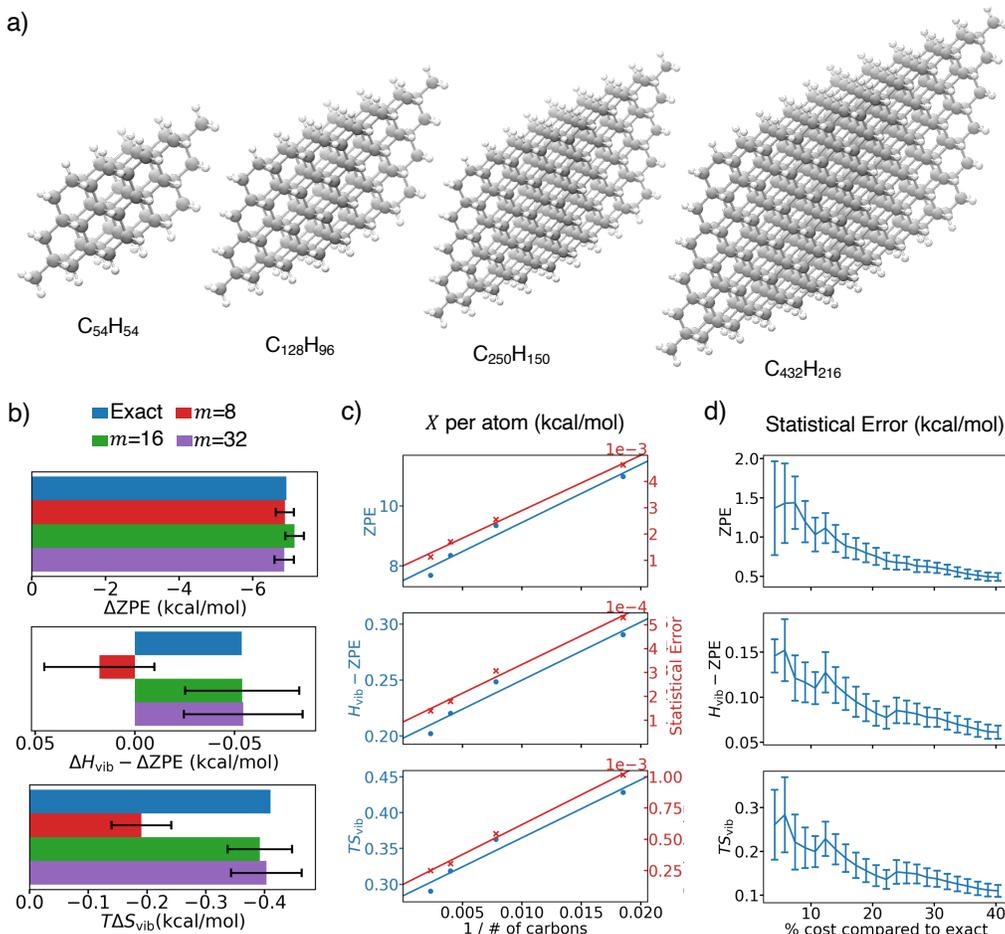
**Fig. 1** **a)** Structures and chemical formulae of the diamond nanocrystals used in the calculations. **b)** $\Delta$ values estimated with different Lanczos quadrature orders ($m$) for the smallest system ($C_{54}H_{54}$) compared to the exact values. Error bars were estimated using the standard error of 50 samples. **c)** Per-atom quantity values and errors as a function of system size for a fixed number of samples (50 samples). The solid lines correspond to linear fits to inverse size. Note that the scale used for the quantities and their errors is different for each of the 3 graphs. **d)** Statistical errors (for $m = 16$ quadrature) in absolute thermodynamic quantities of $C_{432}H_{216}$ as a function of % of computational cost of the exact calculation (error bars denote error of error).

sampling approach described above. We do not include explicit water molecules in the simulation: in principle, these could be included at additional cost, or the desolvation contribution to the free energy can be separately estimated by standard continuum methods [15].

We first study a system which is just small enough that exact results at the xTB level can be obtained at a high computational cost: a cutout ($\sim 1600$ atoms) of the human tankyrase 2 (TNKS2) protein with a bound ligand shown in Figure 2a (see Methods for more information). In Figure 2b we show the thermal contributions to the binding enthalpy, entropy and free energy for the stochastic Lanczos quadrature orders $m = 8, 16, 32$ using xTB. We present a detailed check of the Lanczos convergence for a larger set of $m$ values, and across a set of different systems, in the SI (see
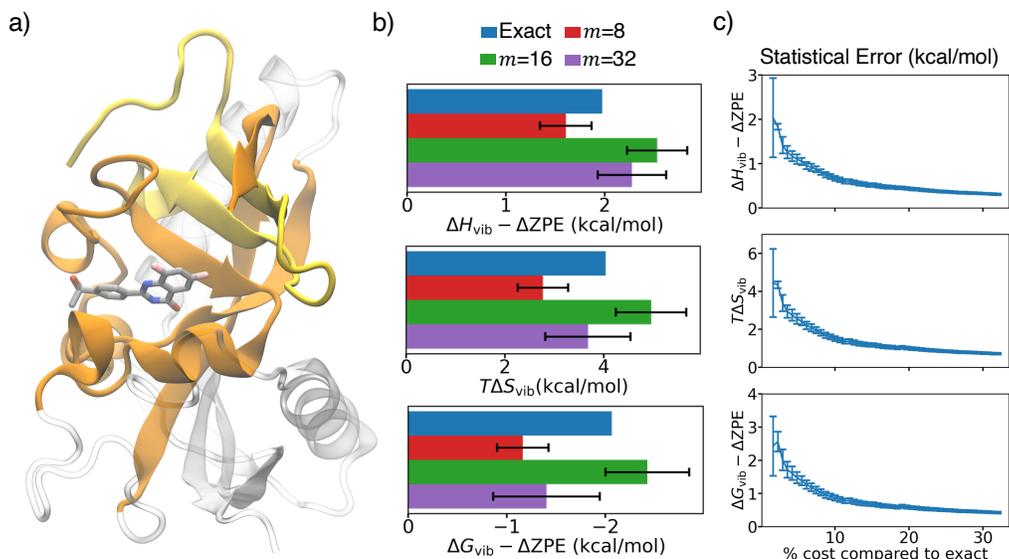
**Fig. 2** Harmonic thermal contributions for the TNKS2 complex computed at the xTB level. **a)** An image of the TNKS2 complex. The truncated part of the protein is shown as transparent, while the remaining two protein chains and the ligand are colored orange, yellow and grey, respectively. Image rendered by VMD[18]. **b)** The thermal enthalpy and entropy, and free energy of binding for the TNKS2 system for varying Lanczos order. The error bars represent $\pm$ one standard error from 100 random samples. **c)** Statistical errors in binding free energy quantities as a function of % of cost of the exact calculation (error bars denote error of error).

Figures 3S and 4S). Following these convergence checks, we choose $m = 16$ for further calculations, where the error due to the Lanczos order is estimated to be less than 1 kcal/mol. Table 1 summarizes the data using up to 100 samples for all rovibrational free energy contributions (the rotational contribution is obtained following Ref. [17]). From comparison to the exact results, the total thermal contribution $\Delta G_{\mathrm{vib}} - \Delta \mathrm{ZPE}$ can be estimated with a statistical error of less than 1 kcal/mol with a cost of roughly 10% of the exact Hessian calculation. The non-thermal contribution $\Delta \mathrm{ZPE}$ (not plotted) has larger statistical error, but can still be estimated to better than 2 kcal/mol with roughly 200 samples, or 67% cost of the exact Hessian calculation.

We next evaluate the thermal quantities at the Kohn-Sham DFT level using the PBE functional for the TNKS2 complex using up to 35 samples, as summarized in Table 1. Interestingly, we find the thermal vibrational contributions at the DFT level to be quite similar to those from xTB. A single sample using our DFT implementation takes roughly two days on 1 node (32 CPU cores), compared to 600 days on 1 node for the exact Hessian calculation.

Finally, we apply our approach to HIV protease bound to a small molecule (JE-2147)[19] (Figure 3a). This system contains over 3000 atoms including hydrogens. The number of gradient calculations required to compute the full Hessian in this case ($\sim 10{,}000$) is so large that the exact computation is expensive even at the level of semi-empirical quantum mechanics, thus we do not compute exact data here. In Table 1 we show our stochastic estimates of the harmonic contributions to the thermodynamic binding quantities using xTB and a quadrature order of $m = 20$. Interestingly, the thermal contributions ($\Delta H_{\mathrm{vib}} - \Delta \mathrm{ZPE}$ and $T\Delta S_{\mathrm{vib}}$) to the free
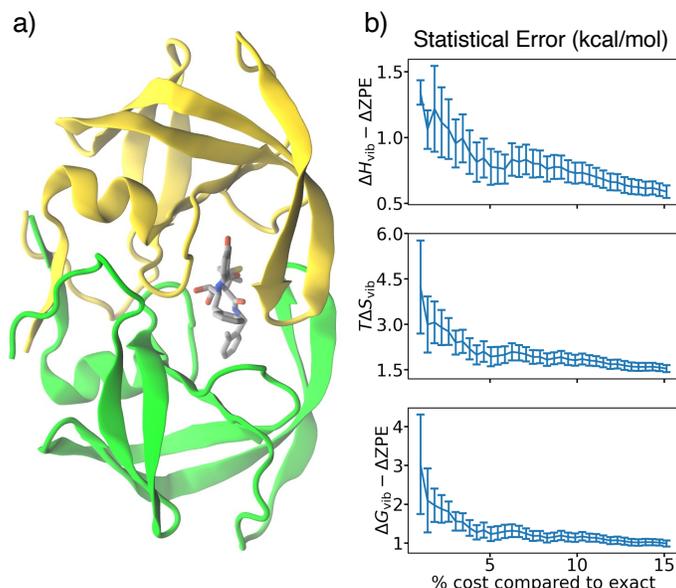
a)

b)



**Fig. 3** Harmonic thermal contributions for the JE-2147-HIV protease complex computed at the xTB level. **a)** An image of the JE-2147-HIV protease complex. The two protein chains and the ligand are shown in yellow, green, and grey respectively. Image rendered by VMD[18]. **b)** Statistical errors in binding free energy quantities as a function of % of cost of the exact calculation (error bars denote error of error).

| Quantity | TNKS (xTB) | TNKS (DFT) | HIV (xTB) | TNKS (expt.) | HIV (expt.) |
|---|---|---|---|---|---|
| $\Delta G_{\mathrm{rot}}$ | 9.57 | 9.61 | 10.43 | | |
| $\Delta H_{\mathrm{vib}} - \Delta \mathrm{ZPE}$ | 2.53±0.30 | 2.89±0.51 | 1.88 ± 0.72 | | |
| $T\Delta S_{\mathrm{vib}}$ | 4.96±0.71 | 5.79±1.11 | 1.69 ± 1.79 | | |
| $\Delta G_{\mathrm{vib}} - \Delta \mathrm{ZPE}$ | -2.43±0.42 | -2.90±0.64 | 0.19 ± 1.13 | | |
| $\Delta G_{\mathrm{bind}}^{\mathrm{tot}}$ | -26.8 | -21.4 | -14.1 | -11.0 | -14.2 |

**Table 1** Contributions to the binding free energy at 298.15 K for the TNKS2 and JE-2147-HIV protease system. The rotational free energies $\Delta G_{\mathrm{rot}}$ are computed at the optimized structures, the thermal enthalpy, entropy and free energy are computed from stochastic sampling. The error estimate for TNKS using xTB is one standard error from 100 random samples (corresponding to 33% of exact cost), 35 samples (corresponding to 11% of exact cost) for TNKS using DFT, and 50 random samples (corresponding to 10% of exact cost) for HIV (xTB). The experimental binding affinity is estimated as either $k_B T \ln \mathrm{IC}_{50}$ or $k_B T \ln K_i$, with $\mathrm{IC}_{50} = 9.2$ nM for TNKS2 from Ref. [20] and $K_i = 41$ pM for HIV protease from Ref. [19]. $\Delta G_{\mathrm{bind}}^{\mathrm{tot}}$ includes the non-thermal contribution from the ZPE and is computed as described in the Methods section. All values are given in kcal/mol.

energy are both similar, small, and of opposite sign, meaning that the total thermal free energy contribution is almost zero. The small size of the thermo-statistical harmonic contributions may be due to the known rigidity of the HIV protease binding pocket, which means that many of the normal modes in the free protein and protein-ligand complex may be very similar. Nonetheless, estimating $\Delta G_{\mathrm{vib}} - \mathrm{ZPE}$ to an accuracy of 1 kcal/mol is clearly feasible within our scheme at roughly 10% of the estimated cost of the exact calculation (Figure 3b).

| Quantity | All Diamonds | TNKS (xTB) | TNKS (DFT) | HIV (xTB) |
|---|---|---|---|---|
| $\Delta H_{\mathrm{vib}} - \Delta\mathrm{ZPE}$ | $<1\%$ | 3% | 3% | 11% |
| $T\Delta S_{\mathrm{vib}}$ | $<1\%$ | 16% | 14% | 74% |
| $\Delta G_{\mathrm{vib}} - \Delta\mathrm{ZPE}$ | $<1\%$ | 6% | 5% | 31% |

**Table 2**  Summary of computational cost to achieve an accuracy of 1 kcal/mol compared to an exact Hessian calculation. The number of samples $n$ required to achieve the desired accuracy is estimated with 50 samples for the diamond systems ($n' = 50$), $n' = 100$ for TNKS2 (xTB), $n' = 35$ for TNKS (DFT), and $n' = 50$ for HIV (xTB). See SI for details about error analysis.

## Discussion

We have presented results which demonstrate the feasibility of computing harmonic contributions to the free energy at the quantum mechanical level for systems of more than a thousand atoms. The cost is greatly reduced from that needed to compute the Hessian of the system. This is particularly true when one is interested in intensive (or "per-atom") quantities, where self-averaging behavior shows that in large systems, we may estimate the quantities at a cost comparable to that of a few energy evaluations. This holds promise in evaluating thermodynamic transitions in materials involving large unit cells, for example, those associated with alloys and disorder. In the case of free energy differences, the correlated sampling technique employed here makes the evaluation of even small thermal free energy differences, as found in protein-ligand complexes, feasible at the level of 1-2 kcal/mol. The estimated speedups for all systems considered, in order to reach a given accuracy, are summarized in Table 2.

An additional advantage of the current approach is that the cost may be continually tuned. This is relevant to new applications, for example in the computational screening for therapeutics [1, 21], where less precise estimates are an acceptable tradeoff for speed. Also, as we see from Table 1, while the computed binding affinity using the harmonic approximation is not always highly accurate in biomolecular systems, the increased facility to obtain harmonic estimates further raises the possibility for new approaches to compute anharmonic contributions to free energies with a variety of techniques, such as the minimum-mining technique [22], which samples multiple minima in an anharmonic potential and combines harmonic contributions from each of them. While applications which require intensive quantities benefit most from the stochastic approach, converging absolute thermal quantities to sufficiently high precision may require improved statistical estimators [23]. In addition, while we have estimated harmonic thermal contributions using quantum mechanical energy functions, the same algorithm accelerates harmonic free energy computation using any energy function, including classical force-fields, and can be combined with other cost reduction techniques, such as partial Hessians.

In summary, the technique presented here suggests that the estimation of harmonic free energy effects at the quantum mechanical level for systems with hundreds or even more than a thousand of atoms need not be considered a future challenge [1], but one which can be begin to be addressed today.

# Methods

## Stochastic Lanczos quadrature

The stochastic Lanczos method is a numerical method which has been employed in different contexts (see e.g. Ref. [24] for an early application in quantum many-body systems). We follow the general mathematical formulation in Ref. [7]. The Lanczos iterations were performed starting from a vector randomly selected from a Rademacher distribution. The Lanczos iterations require the action of the mass-weighted Hessian matrix on this random vector. We compute this matrix-vector product from finite difference gradient calculations:

$$\mathbf{D}\mathbf{v} = \mathbf{M}^{-1/2}\frac{\mathbf{g}(\delta\mathbf{v}) - \mathbf{g}(-\delta\mathbf{v})}{2\delta}. \tag{6}$$

Here, $\mathbf{D}$ is the mass-weighted Hessian matrix, $\mathbf{g}$ is the gradient, $\mathbf{M}$ is the diagonal matrix of masses. The displacement is given by

$$\delta\mathbf{v} = \mathbf{M}^{-1/2}\mathbf{v} \tag{7}$$

where the factor of $\mathbf{M}^{-1/2}$ accounts for the mass-weighting. The value of $\delta$ is chosen based on the norm of the random vector so that the average displacement per atom is 0.0012 Å.

The number of Lanczos iterations $m$ is a parameter of the method; $m$ should be increased until convergence is reached. In chemical systems, the maximum eigenvalue of the Hessian does not scale with the system size (it is the maximum vibrational frequency, for example a C-H stretch), while the minimum eigenvalue is bounded from below by $0$. For many functions of the Hessian, this means that the maximum and minimum eigenvalues also do not scale with system size. Under this assumption, we expect a fixed $m$ to yield a constant relative error in the trace of the function, and furthermore, due to exponential convergence in $m$ for well-behaved functions of the Hessian, $m$ needs to only increase logarithmically with system size for constant absolute error. A numerical study of convergence with $m$ is presented in the SI.

Additionally, we also implemented and tested a Chebyshev fitting method as an alternative to the stochastic Lanczos quadrature. We found that often a higher order Chebyshev fit was required making it a slightly more expensive alternative to Lanczos quadrature.

## Calculations on diamond nanocrystals

Diamond nanocrystals were constructed by creating supercells of the bulk diamond unit cell and then capping with hydrogens. The resulting structure was optimized using the PBE functional[13] and the def2-SV(P) basis set[25]. All DFT calculations were performed with the ORCA program[26, 27]. The structure was also optimized using the second generation extended tight-binding (GFN2-xTB)

method[28] as implemented in the Semiempirical Extended Tight-Binding (xTB) program package[29].

## Calculations on protein-ligand systems

All calculations on protein-ligand systems used the second generation extended tight-binding (GFN2-xTB) method[28] as implemented in the Semiempirical Extended Tight-Binding (xTB) program package[29]. Generalized Born, solvent-accesible area (GBSA) solvation was used to mimic an aqueous environment for all calculations. For the TNKS2 system, we additionally performed the calculations using density functional theory. We used the PBE functional with the GTH-DZV basis and GTH pseudo-potential[30] for PBE. The system was placed in a $45.55\text{Å} \times 41.78\text{Å} \times 37.27\text{Å}$ periodic box, allowing a 5-6Å vacuum around the atoms. A 0.15 Hartree level shift was applied to the virtual orbitals to help the SCF convergence. The Gaussian and Plane Waves method[31, 32] was employed and the plane wave cutoff was 200 Hartree.

The truncated TNKS2 protein was constructed from the the ligands/protein-structure obtained from Ref. [33]. The entire protein was minimized using Amber-Tools, using the Generalized Born implicit, igb=5), the Amber 14 force field[34], and the general AMBER force field (GAFF)[35] for ligands, assigned using Antechamber from AmberTools[36]. Following minimization, truncation and capping of the terminals were carried out using PyMol [37]. Truncation was performed to remove all protein atoms beyond 3-4Å around the ligand. Truncated ends were capped using ACE/NME terminal patches. The ligand bound to the protein is one of the many inhibitors identified in Ref. [38] whose structure is available in the Protein Data Bank [39](PDB: JKN).

The structure of the JE-2147-HIV protease complex was obtained from PDB 1KZK. Hydrogens were added using UCSF Chimera[40] and the structure was optimized first using the GFN-FF force field as implemented in the xTB program package[29] and finally with the GFN2-xTB method[28] ultimately used for harmonic vibrational analysis.

The total binding free energy is estimated as $\Delta G_{\mathrm{rot}} + \Delta H_{\mathrm{vib}} - T\Delta S_{\mathrm{vib}} + \Delta E + \Delta G_{\mathrm{solv}}$, where $\Delta E$ is the single point energy difference and $\Delta G_{\mathrm{solv}}$ is estimated via GBSA within xTB.

**Supplementary information.** Supplementary information on statistical analysis, additional analysis on three transition metal complexes.

## Declarations

- Availability of data and materials: Data is available from the authors upon reasonable request.
- Authors' contributions: AFW and GKC conceived the project. AFW, CL carried out the work. All authors contributed to the writing of the paper.

# References

[1] Grimme, S., Schreiner, P.R.: Computational chemistry: the fate of current methods and future challenges. Angewandte Chemie International Edition **57**(16), 4170–4176 (2018)

[2] Li, H., Jensen, J.H.: Partial Hessian vibrational analysis: The localization of the molecular vibrational energy and entropy. Theoretical Chemistry Accounts **107**(4), 211–219 (2002). https://doi.org/10.1007/s00214-001-0317-7

[3] Woodcock, H.L., Zheng, W., Ghysels, A., Shao, Y., Kong, J., Brooks, B.R.: Vibrational subsystem analysis: A method for probing free energies and correlations in the harmonic limit. Journal of Chemical Physics **129**(21) (2008). https://doi.org/10.1063/1.3013558

[4] Filippone, F., Parrinello, M.: Vibrational analysis from linear response theory. Chemical Physics Letters **345**(1-2), 179–182 (2001). https://doi.org/10.1016/S0009-2614(01)00843-0

[5] Karplus, M., Kushick, J.N.: Method for Estimating the Configurational Entropy of Macromolecules. Macromolecules **14**(2), 325–332 (1981). https://doi.org/10.1021/ma50003a019

[6] Brooks, B.R., Janezic, D., Karplus, M.: Harmonic Analysis of Large Systems. I. Methodology. Journal of Computational Chemistry **16**(12), 1522–1542 (1995)

[7] Ubaru, S., Chen, J., Saad, Y.: Fast estimation of $\mathrm{tr}(f(A))$ via stochastic Lanczos quadrature. SIAM Journal on Matrix Analysis and Applications **38**(4), 1075–1099 (2017)

[8] Hellman, O., Steneteg, P., Abrikosov, I.A., Simak, S.I.: Temperature dependent effective potential method for accurate free energy calculations of solids. Physical Review B - Condensed Matter and Materials Physics **87**(10), 1–8 (2013) arXiv:1303.1145. https://doi.org/10.1103/PhysRevB.87.104111

[9] Errea, I., Calandra, M., Mauri, F.: Anharmonic free energies and phonon dispersions from the stochastic self-consistent harmonic approximation: Application to platinum and palladium hydrides. Physical Review B - Condensed Matter and Materials Physics **89**(6), 1–16 (2014) arXiv:1311.3083. https://doi.org/10.1103/PhysRevB.89.064302

[10] Baer, R., Neuhauser, D., Rabani, E.: Self-averaging stochastic kohn-sham

density-functional theory. Physical Review Letters **111**(10), 1–5 (2013). https://doi.org/10.1103/PhysRevLett.111.106402

[11] Han, I., Malioutov, D., Shin, J.: Large-scale log-determinant computation through stochastic chebyshev expansions. In: International Conference on Machine Learning, pp. 908–917 (2015). PMLR

[12] Kaledin, A.L.: Gradient-based direct normal-mode analysis. The Journal of chemical physics **122**(18), 184106 (2005)

[13] Perdew, J.P., Burke, K., Ernzerhof, M.: Generalized gradient approximation made simple. Physical review letters **77**(18), 3865 (1996)

[14] Grimme, S., Bannwarth, C., Shushkov, P.: A robust and accurate tight-binding quantum chemical method for structures, vibrational frequencies, and noncovalent interactions of large molecular systems parametrized for all spd-block elements (z=1–86). J. Chem. Theory Comput. **13**(5), 1989–2009 (2017). https://doi.org/10.1021/acs.jctc.7b00118

[15] Ehrlich, S., Göller, A.H., Grimme, S.: Towards full quantum-mechanics-based protein–ligand binding affinities. ChemPhysChem **18**(8), 898–905 (2017)

[16] Spicher, S., Grimme, S.: Efficient computation of free energy contributions for association reactions of large molecules. The Journal of Physical Chemistry Letters **11**(16), 6606–6611 (2020)

[17] Grimme, S.: Supramolecular binding thermodynamics by dispersion-corrected density functional theory. Chemistry–A European Journal **18**(32), 9955–9964 (2012)

[18] Humphrey, W., Dalke, A., Schulten, K.: Vmd: visual molecular dynamics. Journal of molecular graphics **14**(1), 33–38 (1996)

[19] Reiling, K.K., Endres, N.F., Dauber, D.S., Craik, C.S., Stroud, R.M.: Anisotropic dynamics of the JE-2147-HIV protease complex: Drug resistance and thermodynamic binding mode examined in a 1.09 Å structure. Biochemistry **41**(14), 4582–4594 (2002). https://doi.org/10.1021/bi011781z

[20] Buchstaller, H.-P., Anlauf, U., Dorsch, D., Kuhn, D., Lehmann, M., Leuthner, B., Musil, D., Radtki, D., Ritzert, C., Rohdich, F., *et al.*: Discovery and optimization of 2-arylquinazolin-4-ones into a potent and selective tankyrase inhibitor modulating wnt pathway activity. Journal of Medicinal Chemistry **62**(17), 7897–7909 (2019)

[21] Mardirossian, N., Wang, Y., Pearlman, D.A., Chan, G.K., Shiozaki, T.: Novel algorithms and high-performance cloud computing enable efficient fully quantum mechanical protein-ligand scoring. arXiv preprint arXiv:2004.08725

(2020)

[22] Chen, W., Gilson, M.K., Webb, S.P., Potter, M.J.: Modeling protein- ligand binding by mining minima. Journal of chemical theory and computation **6**(11), 3540–3557 (2010)

[23] Meyer, R.A., Musco, C., Musco, C., Woodruff, D.P.: Hutch++: Optimal stochastic trace estimation. In: Symposium on Simplicity in Algorithms (SOSA), pp. 142–155 (2021). SIAM

[24] Jaklič, J., Prelovšek, P.: Lanczos method for the calculation of finite-temperature quantities in correlated systems. Physical Review B **49**(7), 5065 (1994)

[25] Weigend, F., Ahlrichs, R.: Balanced basis sets of split valence, triple zeta valence and quadruple zeta valence quality for h to rn: Design and assessment of accuracy. Physical Chemistry Chemical Physics **7**(18), 3297–3305 (2005)

[26] Neese, F.: The ORCA program system. Wiley Interdisciplinary Reviews: Computational Molecular Science **2**(1), 73–78 (2012). https://doi.org/10.1002/wcms.81

[27] Neese, F.: Software update: the ORCA program system, version 4.0. Wiley Interdisciplinary Reviews: Computational Molecular Science **8**(1), 4–9 (2018). https://doi.org/10.1002/wcms.1327

[28] Bannwarth, C., Ehlert, S., Grimme, S.: GFN2-xTB - An Accurate and Broadly Parametrized Self-Consistent Tight-Binding Quantum Chemical Method with Multipole Electrostatics and Density-Dependent Dispersion Contributions. Journal of Chemical Theory and Computation **15**(3), 1652–1671 (2019). https://doi.org/10.1021/acs.jctc.8b01176

[29] Bannwarth, C., Caldeweyher, E., Ehlert, S., Hansen, A., Pracht, P., Seibert, J., Spicher, S., Grimme, S.: Extended tight-binding quantum chemistry methods. Wiley Interdisciplinary Reviews: Computational Molecular Science **11**(2), 1–49 (2021). https://doi.org/10.1002/wcms.1493

[30] Hartwigsen, C., Gœdecker, S., Hutter, J.: Relativistic separable dual-space gaussian pseudopotentials from h to rn. Physical Review B **58**(7), 3641 (1998)

[31] Lippert, B.G., PARRINELLO, J.H., MICHELE: A hybrid gaussian and plane wave density functional scheme. Molecular Physics **92**(3), 477–488 (1997)

[32] McClain, J., Sun, Q., Chan, G.K.-L., Berkelbach, T.C.: Gaussian-based coupled-cluster theory for the ground-state and band structure of solids. Journal of chemical theory and computation **13**(3), 1209–1218 (2017)

[33] Schindler, C.E.M., Baumann, H., Blum, A., Böse, D., Buchstaller, H.P., Burgdorf, L., Cappel, D., Chekler, E., Czodrowski, P., Dorsch, D., Eguida, M.K.I., Follows, B., Fuchß, T., Grädler, U., Gunera, J., Johnson, T., Lebrun, C.J., Karra, S., Klein, M., Knehans, T., Koetzner, L., Krier, M., Leiendecker, M., Leuthner, B., Li, L., Mochalkin, I., Musil, D., Neagu, C., Rippmann, F., Schiemann, K., Schulz, R., Steinbrecher, T., Tanzer, E.M., Lopez, A.U., Follis, A.V., Wegener, A., Kuhn, D.: Large-scale assessment of binding free energy calculations in active drug discovery projects. Journal of Chemical Information and Modeling **60**(11), 5457–5474 (2020). https://doi.org/10.1021/acs.jcim.0c00900

[34] Maier, J.A., Martinez, C., Kasavajhala, K., Wickstrom, L., Hauser, K.E., Simmerling, C.: ff14SB: Improving the Accuracy of Protein Side Chain and Backbone Parameters from ff99SB. Journal of Chemical Theory and Computation **11**(8), 3696–3713 (2015). https://doi.org/10.1021/acs.jctc.5b00255

[35] Wang, J., Wolf, R.M., Caldwell, J.W., Kollman, P.A., Case, D.A.: Development and testing of a general Amber force field. Journal of Computational Chemistry **25**(9), 1157–1174 (2004). https://doi.org/10.1002/jcc.20035

[36] D.A. Case, H.M. Aktulga, K. Belfon, I.Y. Ben-Shalom, S.R. Brozell, D.S. Cerutti, T.E. Cheatham, III, V.W.D. Cruzeiro, T.A. Darden, R.E. Duke, G. Giambasu, M.K. Gilson, H. Gohlke, A.W. Goetz, R. Harris, S. Izadi, S.A. Izmailov, C. Jin, K. Kasavajhala, M.C., Kollman, P.A.: Amber 2021, Univeristy of California, San Francisco (2021). https://ambermd.org/doc12/Amber21.pdf

[37] Schrödinger, L.: The PyMOL molecular graphics system, version 1.8. November (2015)

[38] Waaler, J., Leenders, R.G.G., Sowa, S.T., Alam Brinch, S., Lycke, M., Nieczypor, P., Aertssen, S., Murthy, S., Galera-Prat, A., Damen, E., Wegert, A., Nazaré, M., Lehtiö, L., Krauss, S.: Preclinical Lead Optimization of a 1,2,4-Triazole Based Tankyrase Inhibitor. Journal of Medicinal Chemistry **63**(13), 6834–6846 (2020). https://doi.org/10.1021/acs.jmedchem.0c00208

[39] Berman, H.M., Battistuz, T., Bhat, T.N., Bluhm, W.F., Bourne, P.E., Burkhardt, K., Feng, Z., Gilliland, G.L., Iype, L., Jain, S., Fagan, P., Marvin, J., Padilla, D., Ravichandran, V., Schneider, B., Thanki, N., Weissig, H., Westbrook, J.D., Zardecki, C.: The protein data bank. Acta Crystallographica Section D: Biological Crystallography **58**(6 I), 899–907 (2002). https://doi.org/10.1107/S0907444902003451

[40] Pettersen, E.F., Goddard, T.D., Huang, C.C., Couch, G.S., Greenblatt, D.M., Meng, E.C., Ferrin, T.E.: UCSF Chimera - A visualization system for exploratory research and analysis. Journal of Computational Chemistry **25**(13), 1605–1612 (2004). https://doi.org/10.1002/jcc.20084

# Supplementary Information for "Quantum harmonic free energies for biomolecules and nanomaterials"

Alec F. White[1,2], Chenghan Li[1], Xing Zhang[1] and Garnet Kin-Lic Chan[1]

[1]Division of Chemistry and Chemical Engineering, California Institute of Technology, 1200 East California Boulevard, Pasadena, 91125, CA, United States of America.
[2]Present Address: Quantum Simulation Technologies, Inc., Cambridge, MA 02139, United States of America.
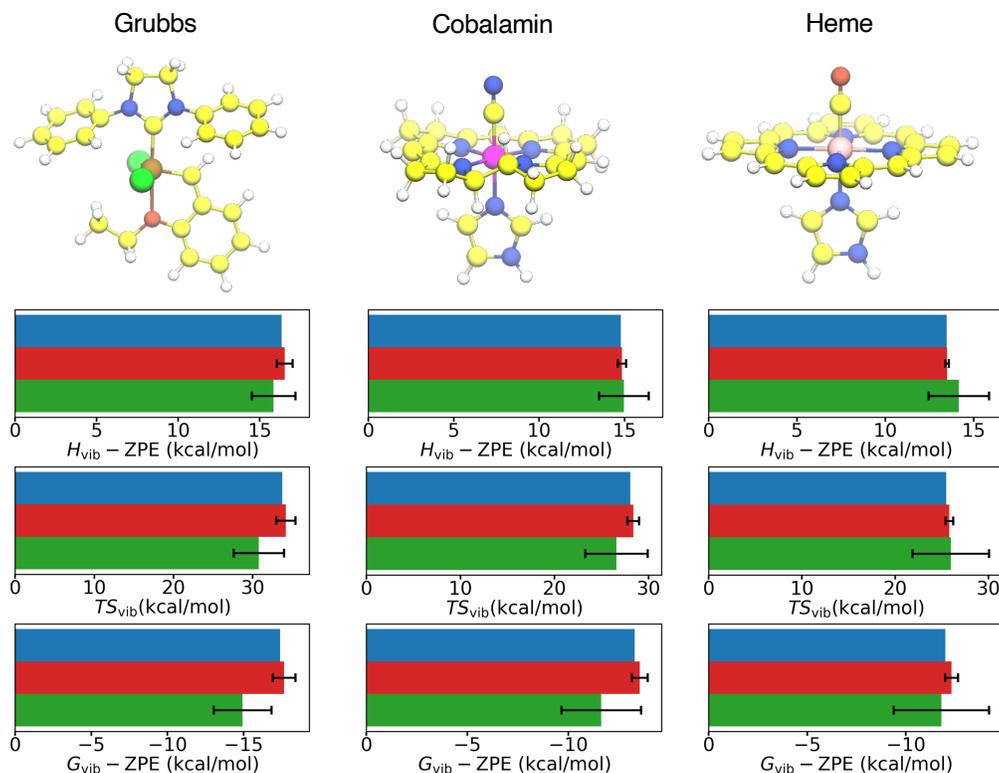
## 1 Error estimation

For each random vector $\mathbf{v}_l$, $s_l = M \times \mathbf{v}_l^T f(\mathbf{D})\mathbf{v}_l$ forms an unbiased estimator of $\text{Tr}(f(\mathbf{D}))$, and so does the average of $n$ samples, $\bar{s}_n = \frac{1}{n}\sum_l^n s_l$. Since each random vector is drawn independently, it is well-known that the sample variance $S^2$ (defined below) is an unbiased estimator of $\text{Var}(\bar{s}_n)$, i.e., $\mathbb{E}[(\bar{s}_n - \mathbb{E}[\bar{s}_n])^2]$, and the stochastic error in the $n$-sample estimator $\bar{s}_n$ can be estimated as follows

$$S_n = \sqrt{\frac{1}{n(n-1)}\sum_l^n (s_l - \bar{s}_n)^2} \tag{1S}$$

Such an $S_n$ is sometimes referred to as the standard error, and we follow this convention in our main text. To estimate the error of $S_n$, we would like to estimate $\text{Var}(S_n^2) = \frac{1}{n}(\mu_4 - \frac{n-3}{n-1}\mu_2^2)$ from our finite samples, where $\mu_p = \mathbb{E}[(s_l - \mathbb{E}[s_l])^p]$. We use unbiased estimators for $\mu_4$ and $\mu_2^2$ (sometimes known as (generalized) $h$-statistics)

$$h_4 = \frac{(n^3 - 2n^2 + 3n)m_4 - n(6n-9)m_2^2}{(n-1)(n-2)(n-3)} \tag{2S}$$

1

Supplementary Fig. 1: Thermal quantities at the DFT-level for a Grubbs catalyst, cobalamin, and heme. The blue bars indicate the exact DFT values, the red bars represent the average over 10 samples from correlated sampling using xTB as the low-level method, and the green bars are the average from 10 samples of direct sampling with DFT. The error bars are estimates for a 2-sample estimator computed from the total 10 samples. For these relatively small systems, 2 samples correspond to roughly 20% of the cost of computing exact numerical Hessian.

$$h_{2,2} = \frac{-(n^2 - n)m_4 + n(n^2 - 3n + 3)m_2^2}{(n-1)(n-2)(n-3)} \tag{3S}$$

where $m_p$ is the $p$-th central moment of the samples, such that $\frac{1}{n}\left(h_4 - \frac{n-3}{n-1}h_{2,2}\right)$ forms an unbiased estimator for $\mathrm{Var}(S_n^2)$, and the error of error is computed from its square root.

We next consider error estimation for $\bar{s}_n$ from $n'$ samples where $n' > n$. A special case is to estimate the error of a single-sample estimator $\bar{s}_1 = s_l$ from $n'$ samples. Such an error is just the variance of the underlying distribution of $s_l$, i.e., $\mathbb{E}[(s_l - \mathbb{E}[s_l])^2]$, and it is straightforward to show that $n'S_{n'}^2$ forms an unbiased estimator. Similarly, one can show that $n'S_{n'}^2/n$ is an unbiased estimator from $n'$ samples for $\mathrm{Var}(\bar{s}_n)$.

## 2 Performance on transition-metal complexes

Transition metal complexes provide an example of systems where empirical force fields are often inadequate and quantum mechanical energy functions are especially

| System | | ZPE | $H_{vib} - ZPE$ | $TS_{vib}$ |
|---|---|---|---|---|
| $C_{54}H_{54}$ | Exact | 599.8436 | 15.742 | 23.5286 |
| | Sampled | $602\pm4$ | $16.0\pm0.3$ | $24.0\pm0.6$ |
| $C_{128}H_{96}$ | Exact | 1206.4788 | 31.9561 | 47.1107 |
| | Sampled | $1209\pm5$ | $32.0\pm0.4$ | $47.2\pm0.6$ |

Supplementary Table 1: Absolute quantities of diamonds showing self-averaging.

valuable. We considered three transition-metal complexes, namely a Grubbs catalyst, a cobalamin model, and a heme model, as challenging cases. The initial structure of the Grubbs catalyst was taken from Ref. [1], optimized with both xTB and B3LYP[2] functional theory. The 6-311G basis set[3, 4] was used for main group elements while the LanL2DZ basis[5] was used for the transition metal Ru. The stochastic sampling was performed using the same level of theory. The cobalamin structure was taken from Ref. [6], and BP86[7, 8]/6-31G(d)[9–11] was used for geometry optimization and sampling. The heme structure was taken from Ref. [12], and we used the B3LYP/6-311G(d)[3] level of theory. The DFT methods were chosen to be similar to the ones employed in the references from which the initial structures were obtained. In all the calculations, a Lanczos order of 16 was used, which we see to be sufficient to reproduce the exact results to within the small statistical error bars.

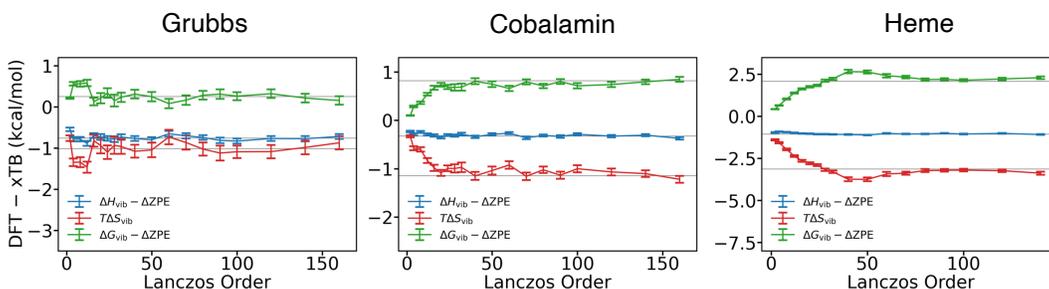# 3 Self-averaging in diamond nanocrystals

In order to confirm that the observed small statistical error in large diamond crystals is the result of "self-averaging", instead of due to under-sampling of the larger coordinate space, we show in Supplementary Table 1 the exact absolute free energy of the two diamond systems, computed from the exact xTB Hessian and from our method using 50 samples (numbers in kcal/mol). We see that the exact values all lie within one standard error of the sampled values, validating the error bars at both system sizes, despite the fixed number of samples as a function of system size. Moving from $C_{54}H_{54}$ to $C_{128}H_{96}$, the system size more than doubles, but the statistical error increases by much less (and in fact stays nearly constant), which is the self-averaging effect referred to in the text.

# 4 Convergence with fitting order

Here, we check the convergence with respect to the fitting order $m$. To obtain a basic understanding, we first examine the accuracy of Chebyshev polynomial fitting for the thermodynamic quantities, which can be assessed without using any stochastic sampling. We do so by examining the accuracy of the Chebyshev expansion over a range of frequencies (see Supplementary Figure 2). The Chebyshev error for a specific system can be estimated by averaging this fitting error over the frequency domain, weighted by the density of the system's vibrational states. The lowest vibrational frequency is $0 \text{ cm}^{-1}$, while the highest frequency of a chemical system is usually the bond vibration between a hydrogen and a heavy atom (typically $< 4000 \text{ cm}^{-1}$). Thus the only system dependence comes from the density of states. This is in contrast to the case of Chebyshev fitting in electronic structure, where the maximum frequency

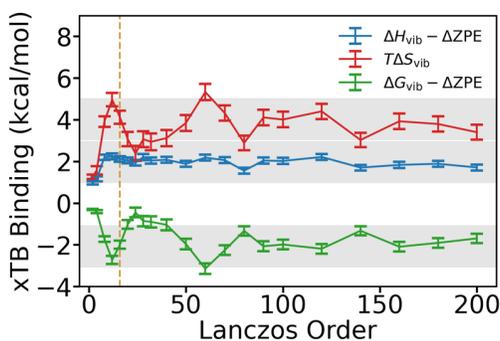Supplementary Fig. 2: Chebyshev fitting to thermodynamic quantities at 298.15 K as a function of fitting order $m$.



Supplementary Fig. 3: Convergence w.r.t. Lanczos order $m$ for transition metal complexes. Grey lines indicate the exact values. Stochastic errors are estimated from 500 samples.

usually grows with system size. We see in both plots that the maximum pointwise deviation at $m = 16$ is less than 1 kcal/mol over the range plotted (although we note that the entropy diverges at $0$ cm$^{-1}$ where the density of states also vanishes). The above convergence check cannot be directly carried out for stochastic Lanczos quadrature due to the need to specify some initial stochastic vector. We have therefore checked the convergence of the Lanczos method as a function of $m$ for several systems by brute-force stochastic sampling, as shown in Supplementary Figures 3 and 4. A Lanczos order of $m = 16$ is generally seen to be sufficient to obtain a systematic error of 1 kcal/mol or less.

# References

[1] Albalawi, M.O., Falivene, L., Jedidi, A., Osman, O.I., Elroby, S.A., Cavallo, L.: Influence of the anionic ligands on properties and reactivity of hoveyda-grubbs catalysts. Molecular Catalysis **509**, 111612 (2021)

[2] Stephens, P.J., Devlin, F.J., Chabalowski, C.F., Frisch, M.J.: Ab initio calculation of vibrational absorption and circular dichroism spectra using density functional force fields. The Journal of physical chemistry **98**(45), 11623–11627 (1994)

Supplementary Fig. 4: Convergence w.r.t. Lanczos order $m$ for TNKS ligand binding. Grey shaded area indicates $\pm$ 1 kcal/mol around the exact values. The vertical dashed line indicates $m = 16$. Stochastic errors are estimated from 200 samples.

[3] Krishnan, R., Binkley, J.S., Seeger, R., Pople, J.A.: Self-consistent molecular orbital methods. xx. a basis set for correlated wave functions. The Journal of chemical physics **72**(1), 650–654 (1980)

[4] McLean, A., Chandler, G.: Contracted gaussian basis sets for molecular calculations. i. second row atoms, z= 11–18. The Journal of chemical physics **72**(10), 5639–5648 (1980)

[5] Hay, P.J., Wadt, W.R.: Ab initio effective core potentials for molecular calculations. potentials for k to au including the outermost core orbitals. The Journal of chemical physics **82**(1), 299–310 (1985)

[6] Kornobis, K., Kumar, N., Wong, B.M., Lodowski, P., Jaworska, M., Andruniów, T., Ruud, K., Kozlowski, P.M.: Electronically excited states of vitamin b12: benchmark calculations including time-dependent density functional theory and correlated ab initio methods. The Journal of Physical Chemistry A **115**(7), 1280–1292 (2011)

[7] Becke, A.D.: Density-functional exchange-energy approximation with correct asymptotic behavior. Physical review A **38**(6), 3098 (1988)

[8] Perdew, J.P.: Density-functional approximation for the correlation energy of the inhomogeneous electron gas. Physical Review B **33**(12), 8822 (1986)

[9] Hehre, W.J., Ditchfield, R., Pople, J.A.: Self—consistent molecular orbital methods. xii. further extensions of gaussian—type basis sets for use in molecular orbital studies of organic molecules. The Journal of Chemical Physics **56**(5), 2257–2261 (1972)

[10] Hariharan, P.C., Pople, J.A.: The influence of polarization functions on molecular orbital hydrogenation energies. Theoretica chimica acta **28**(3), 213–222 (1973)

[11] Rassolov, V.A., Pople, J.A., Ratner, M.A., Windus, T.L.: 6-31g* basis set for atoms k through zn. The Journal of chemical physics **109**(4), 1223–1229 (1998)

[12] Harvey, J.N.: Dft computation of the intrinsic barrier to co geminate recombination with heme compounds. Journal of the American Chemical Society **122**(49), 12401–12402 (2000)