

## MIT Open Access Articles

*Exploiting graphical processing units to enable quantum chemistry calculation of large solvated molecules with conductor-like polarizable continuum models*

The MIT Faculty has made this article openly available. **Please share** how this access benefits you. Your story matters.

**Citation:** Liu, Fang et al. "Exploiting graphical processing units to enable quantum chemistry calculation of large solvated molecules with conductor-like polarizable continuum models." International Journal of Quantum Chemistry 119, 1 (January 2019): e25760 © 2018 Wiley

**As Published:** <http://dx.doi.org/10.1002/qua.25760>

**Publisher:** Wiley

**Persistent URL:** <https://hdl.handle.net/1721.1/123830>

**Version:** Author's final manuscript: final author's manuscript post peer review, without publisher's formatting or copy editing

**Terms of use:** Creative Commons Attribution-Noncommercial-Share Alike



# Exploiting Graphical Processing Units to Enable Quantum Chemistry Calculation of Large Solvated Molecules with Polarizable Continuum Models

Fang Liu,<sup>1,2,3</sup> David M. Sanchez,<sup>1,2</sup> Heather J. Kulik,<sup>3</sup> Todd J. Martínez<sup>\*1,2</sup>

<sup>1</sup>*Department of Chemistry and The PULSE Institute, Stanford University, Stanford, CA 94305*

<sup>2</sup>*SLAC National Accelerator Laboratory, Menlo Park, CA 94025*

<sup>3</sup>*Department of Chemical Engineering, Massachusetts Institute of Technology, Cambridge, Massachusetts 02139, USA*

**Abstract:** The conductor-like polarizable continuum model (C-PCM) with switching/Gaussian smooth discretization is a widely used implicit solvation model in quantum chemistry. We have previously implemented C-PCM solvation for Hartree-Fock (HF) and density functional theory (DFT) on graphical processing units (GPUs), enabling the quantum mechanical treatment of large solvated biomolecules. Here, we first propose a GPU-based algorithm for the PCM conjugate gradient linear solver that greatly improves the performance for very large molecules. The overhead for PCM related evaluations now consumes less than 15% of the total runtime for DFT calculations on large molecules. Second, we demonstrate that our algorithms tailored for ground state C-PCM are transferrable to excited state properties. Using a single GPU, our method evaluates the analytic gradient of the linear response PCM time-dependent density functional theory (TDDFT) energy up to 80X faster than a conventional central processing unit (CPU)-based implementation. In addition, our C-PCM algorithms are transferable to other methods that require electrostatic potential evaluations. For example, we achieve speedups of up to 130x for restricted electrostatic potential (RESP) based atomic charge evaluations, compared to CPU-based codes. We also summarize and compare the different PCM cavity discretization schemes used in some popular quantum chemistry packages as a reference for both users and developers.

## 1. Introduction

The evolution of computational hardware and algorithms has enabled quantum mechanical (QM) studies for large molecules, including not only static properties but also dynamic processes. This has led to fully QM calculations of biomolecular dynamics in solvated environments.<sup>1</sup> Implicit solvent models based on a dielectric continuum approximation are an attractive conceptual framework to describe solvent effects within a QM approach due to their high efficiency.<sup>2,3</sup> Methods such as the polarizable continuum model<sup>4</sup> (PCM) and its variants such as conductor-like models (COSMO,<sup>5-7</sup> C-PCM<sup>8</sup>, also known as GCOSMO<sup>9</sup>, and IEF-PCM<sup>10-12</sup>) are the most popular and accurate of these implicit solvent models in the context of *ab initio* quantum chemistry.

In early implementations of PCM methods, geometry optimization and molecular dynamics were troublesome, as most models led to discontinuous gradients of the solvation energy. These discontinuities were caused by the discretization of the cavity surface: apparent surface charges (ASCs) at the boundary of two adjacent atoms were abruptly switched on and off (exposed or buried in the cavity surface), making the solvated potential energy surface non-smooth. This problem has recently been overcome by using improved discretization schemes and Gaussian smearing of the ASCs.<sup>13-15</sup>

Unfortunately, the application of PCM in QM calculations of large-scale biomolecular systems can still be limited by CPU bottlenecks.<sup>16</sup> On the one hand, calculation of the solvent-solute Coulomb interaction involves one-electron integrals (OEI), with effort scaling as  $O(MN^2)$ , where  $M$  and  $N$  denote the number of ASCs and basis functions, respectively. For small proteins with hundreds of atoms, the typical values of  $M$  and  $N$  are of the order of  $10^4$  and  $10^3$ , respectively. On the other hand, the most straightforward algorithm (direct inversion) for solving the linear equations for PCM requires  $O(M^3)$  time and  $O(M^2)$  memory. Many efforts have aimed to make PCM methods more efficient. Fast summation techniques like the fast multipole method<sup>17, 18</sup> (FMM) were introduced to PCM<sup>19</sup> to reduce the cost of evaluating Coulomb interactions. Recently, an elegant domain decomposition strategy (ddCOSMO)<sup>20-24</sup> has been proposed as a linear scaling approach to COSMO and implemented with semiempirical methods.

Recently, we presented a completely different strategy to solve the C-PCM problem for large molecules.<sup>25</sup> We exploited the high parallelizability of GPUs to accelerate the evaluation of PCM-related OEs. The linear equations were solved with a highly efficient Randomized Block-Jacobi (RBJ) preconditioner for conjugate gradient (CG),<sup>26</sup> and we introduced a dynamic CG convergence threshold to greatly reduce the number of CG iterations without affecting the final accuracy. Since the publication of that work,<sup>25</sup> we now also build the matrix-vectors products on-the-fly on the GPUs, greatly reducing the memory requirement for large molecules. With these techniques, we are able to perform C-PCM HF<sup>27</sup> or DFT<sup>28</sup> calculations of energies and gradients for any molecule that is tractable in isolation. The GPU-based PCM-related OEI builds and linear solver techniques are the foundation for accelerating PCM calculations in more advanced quantum chemistry methods, such as TDDFT.<sup>29-32</sup> The algorithms are also transferrable to other methods that require evaluation of the electrostatic potential (ESP), as is often used for determining atomic charges.<sup>33-37</sup>

The article is organized as follows: 1) We briefly review the basic formalism of C-PCM, summarizing the different cavity discretization schemes implemented in various quantum chemistry packages. 2) We then explain the necessity of accelerating C-PCM quantum chemistry calculations with GPUs, and briefly review the strategies we used. 3) We introduce recent improvements to the linear solver for very large molecules. 4) Then we demonstrate the transferability of our algorithms to other quantum chemistry methods and benchmark their performance. 5) Finally, we conclude with comments about future extensions of GPU-accelerated PCM.

## 2. Basic Formalism

### 2.1 Conductor-like Polarizable Continuum Model

Prior to explaining the basics of the C-PCM method, it is useful to review key notation used in this article: atomic orbitals (AOs) will be indexed with  $\mu, \nu, \kappa, \gamma$  and  $\tau$ . Molecular orbitals (MOs) will be indexed with  $i, j$  for occupied orbitals, and  $a, b$  for virtual orbitals. Primitive Gaussian basis functions will be indexed with  $\chi_p, \chi_q$ . We will use brackets to denote integrals over primitive basis functions and parentheses to denote

integrals for contracted basis functions. Energies will be differentiated with respect to an external perturbation (e.g., the nuclear coordinate or an external field) denoted as  $\xi$ .

The basic idea of C-PCM is shown schematically in Figure 1. The reaction potential generated by the solute charge distribution is described in terms of an apparent charge distribution spread over the solute cavity surface. Outside the surface is a dielectric continuum with dielectric constant  $\epsilon_2$ , whereas the cavity containing the solute has unit dielectric constant. For numerical convenience, the cavity boundary is usually discretized into  $M$  surface elements (tesserae). Each of these tesserae has an ASC,  $q_k$ , which describes the electric field of the polarized continuum. The values,  $\mathbf{q}$ , of the ASCs are determined with a set of linear equations:<sup>2, 14</sup>

$$\mathbf{Kq} = \mathbf{RV} \quad (0)$$

where  $\mathbf{K}$  and  $\mathbf{R}$  depend on the particular PCM method, and represent geometric information concerning the cavity and the dielectric screening, respectively, and  $\mathbf{V}$  is a vector that represents the solute electrostatic potential at each tesserae. Then the electrostatic component of the solvation free energy is given by the interaction between the polarization charges and solute, in addition to the self-energy of the surface charges. In the context of a self-consistent field (SCF) calculation, it is given as

$$\Delta G_{\text{solv}} = \frac{1}{2} \mathbf{V}^\dagger \mathbf{q} \quad (0)$$

In a PCM calculation of a QM solute molecule, two quantities must be determined: the polarization charges  $\mathbf{q}$  and the solute electron density matrix  $\mathbf{P}$ . In the simplest case, assuming closed shell HF or KS calculations, the solute density matrix is computed as:

$$\mathbf{F}^{\text{solvated}} \mathbf{C} = \mathbf{S} \mathbf{C} \epsilon \quad (0)$$

$$P_{\mu\nu} = 2 \sum_{i=1}^{\text{occ}} C_{\mu i} C_{\nu i}^* \quad (0)$$

where  $\mathbf{F}^{\text{solvated}}$  is the Fock matrix of the solvated system,  $\mathbf{C}$  are the MO coefficients, and  $\mathbf{S}$  is the AO overlap matrix.

The solvent contributes to  $\mathbf{F}^{\text{solvated}}$  through terms containing cavity charges  $\mathbf{q}$

$$\mathbf{F}^{\text{solvated}} = \mathbf{F}^0 + \Delta\mathbf{F}^S(\mathbf{q}) \quad (0)$$

where  $\mathbf{F}^0$  is the gas phase Fock matrix, and  $\Delta\mathbf{F}^S(\mathbf{q})$  is given as:

$$\Delta\mathbf{F}^S(\mathbf{q}) = \sum_{k=1}^M q_k L_{\mu\nu}^k \quad (0)$$

$$L_{\mu\nu}^k = -\int \hat{J}_k^{\text{screened}} \phi_{\mu}(\vec{r}) \phi_{\nu}(\vec{r}) d\mathbf{r} \quad (0)$$

where  $L_{\mu\nu}^k$  represents the electrostatic potential of the basis pair,  $\phi_{\mu}(r)\phi_{\nu}(r)$  at tessera  $k$ , and the exact form of the operator  $\hat{J}_k^{\text{screened}}$  depends on the cavity discretization scheme. For example, in the crudest treatment, when the ASCs are represented by point charges,

$$\hat{J}_k^{\text{screened}} = \frac{1}{|\mathbf{r} - \mathbf{r}_k|} \quad (0)$$

and  $L_{\mu\nu}^k$  has the same form as the usual nuclear-electron attraction integrals.<sup>38</sup>

To determine  $\mathbf{q}$ , one can construct the ESP,  $V_k$ , for the  $k$ th tessera from contributions of the nuclei and solute electron density:

$$V_k = V_k^{\text{nuc}} + V_k^e \quad (0)$$

$$V_k^{\text{nuc}} = \sum_{J=1}^{n_A} Z_J \int \hat{J}_k^{\text{screened}} \delta(\mathbf{r} - \mathbf{R}_J) d\mathbf{r} \quad (0)$$

$$V_k^e = \sum_{\mu\nu} P_{\mu\nu} L_{\mu\nu}^k \quad (0)$$

where  $V^{\text{nuc}}$  and  $V^e$  are the ESP contributions from the solute nuclei and electrons, respectively;  $\mathbf{R}_J$  is the location of the  $J$ th nucleus with atomic radius  $R_J$ . The construction

of  $\mathbf{V}^e$  is very similar to  $\Delta\mathbf{F}^s(\mathbf{q})$ , but the integrals ( $L_{\mu\nu}^k$ ) are contracted over basis function pairs ( $\mu$  and  $\nu$ ) rather than tesserae ( $k$ ).

Therefore, in a PCM-SCF iteration, the computationally intensive steps include three parts:

- (1) constructing  $\mathbf{V}^e$ , i.e. Eq. (0);
- (2) solving the linear system in Eq. (0);
- (3) constructing  $\Delta\mathbf{F}^s$ , i.e. Eq. (0)

We will revisit how to reduce the computational cost in these steps in Section 2.3 and Section 4.

## 2.2 Comparison of Cavity Surface Discretization Schemes

There exist many versions of PCM in different quantum chemistry packages that give slightly different results, which can sometimes cause confusion. Here we briefly summarize the similarities and differences of these implementations.

The first difference lies in the integral equations, as reflected in the detailed formulas for both the  $\mathbf{K}$  and  $\mathbf{R}$  matrices in Eq. 1. This distinguishes the conductor-like models from other PCM models such as SS(V)PE<sup>39</sup> or IEF-PCM.<sup>10-12</sup> For further detail, see the detailed discussion by Lange and Herbert.<sup>14</sup>

The second difference lies in the cavity surface discretization scheme. Even with the same model, e.g. C-PCM, there still exist many variations in different quantum chemistry packages that are distinguished only by the discretization scheme, as embodied by the  $\mathbf{K}$  matrix in Eq. (0). To define a cavity surface discretization scheme, two elements are needed: cavity type (defining the shape of the molecular surface) and tessellation scheme (discretization of the molecular surface).

**Cavity type:** Most cavity types fall into one of the following three categories: (1) the van de Waals cavity; (2) the solvent accessible surface (SAS);<sup>5</sup> or (3) the solvent excluded surface (SES or Connolly Surface<sup>40</sup>) which are shown schematically in Figure 2. It is worth noting that in the literature, the same terminology may refer to different types of cavities. For instance, SAS was originally defined as shown in Figure 2d.<sup>5</sup> The fused vdW surface (Figure 2b, scaling factor  $\alpha=1$ ) is first expanded by the solvent radius (Figure 2c) and then contracted by a distance  $\delta^{\text{SC}}$ , leading to a cavity

with some creases without tesserae at the boundary between atoms. However, in more recent publications,<sup>13,14</sup> SAS refers to Figure 2b, a fused vdW surface with expanded atomic radii ( $\alpha > 1$ ). The SES removes cusps at boundaries by introducing “probe spheres” (see Figures 2e-f). In our implementation, we always employ the modern definition of SAS (Figure 2b).

**Tessellation scheme:** There are more ways to tessellate than types of cavity construction. In total, there are two important aspects regarding tessellation. 1) How is the surface charge distribution on each tessera represented? 2) How is the molecular surface discretized into small segments? For the first aspect, one could represent the charge distribution as a point charge or a Gaussian distribution centered on the tessera. This choice, combined with the techniques defining exposure and burial of the tesserae, determines the continuity of the analytic gradient of the solvation energy. For the second aspect, the discretization methods can be grouped into two categories: polyhedron based and Lebedev grid<sup>41</sup> based. The original COSMO model by Klamt generates the grid points on a sphere from iterative refinement of triangles starting from an icosahedron.<sup>5</sup> Although developed independently, GEPOL<sup>42-45</sup> has a similar procedure for polyhedron refinement, and is still widely used in modern implementations of C-PCM methods in popular electronic structure packages, such as GAMESS-US,<sup>46, 47</sup> GAUSSIAN,<sup>48</sup> and ORCA.<sup>49</sup> The Lebedev grid, which is often used for numerical quadrature in DFT,<sup>50, 51</sup> was introduced to solvent surface discretization by Karplus and coworkers.<sup>13</sup> The Lebedev grid benefits from its advantages in numerical quadrature stability and ease of construction compared to the polyhedron method. It is used in more recent implementations of C-PCM, for example, in Q-CHEM<sup>14, 52, 53</sup> and TERACHEM.<sup>26-28,34,50</sup> The choices of charge representation and surface discretization are independent and can be combined in different ways. For example, a surface discretized with GEPOL can either be represented by point charges, as implemented in GAMESS-US,<sup>42-45, 54</sup> or by Gaussian-smearred distributions, as implemented in GAUSSIAN.<sup>15</sup> Similarly, the Lebedev grid can be either used in the SWIG/ISWIG scheme to obtain continuous energy gradients, or it can be used with point charges associated with fixed tesserae area. Both of these methods are implemented in Q-CHEM and TERACHEM. A summary of different cavity discretization schemes used in some popular quantum chemistry packages is provided in Table 1.

In the following sections, we focus on the SWIG/ISWIG implementation of C-PCM, which benefits from a simple discretization scheme and a smooth analytical gradient of solvation energy. In this context,  $\mathbf{R}$  in Eq. (0) is a scalar

$$\mathbf{R} = -\frac{\varepsilon - 1}{\varepsilon} \quad (0)$$

and the  $\mathbf{K}$  matrix is given as:

$$K_{kl} = \frac{\text{erf}(\zeta'_{kl} |\vec{r}_k - \vec{r}_l|)}{|\vec{r}_k - \vec{r}_l|} \quad \forall k \neq l \quad (0)$$

$$K_{kk} = \frac{\zeta_k}{\sqrt{2\pi}} \mathbf{S}_k^{-1} \quad (0)$$

where  $\vec{r}_k$  is the location of the  $k$ th Lebedev point,  $\zeta_k = \frac{\zeta}{R_{J(k)} \sqrt{w_k}}$  is the Gaussian

exponent for the  $k$ th point where  $\zeta$  is an optimized exponent for the specific Lebedev quadrature level being used (as tabulated<sup>13</sup> by York and Karplus),  $w_k$  is the Lebedev quadrature weight for the  $k$ th point, and  $R_{J(k)}$  is the radius of the atom ( $J$ ) associated with the  $k$ th point. Finally,  $\zeta'_{kl} = \frac{\zeta_k \zeta_l}{\sqrt{\zeta_k^2 + \zeta_l^2}}$  is the combined exponent that emerges when

describing the Coulomb interaction of two Gaussian charge distributions. The matrix  $\mathbf{S}_k$  in Eq. (0) is the switching function which smoothes the boundary of adjacent vdW spheres corresponding to different atoms (and thus makes the solvation energy continuous), and is a function of only the positions of the PCM grid points and solute nuclei. Using Gaussian distributions to represent the ASC, the operator  $\hat{J}_k^{\text{screened}}$  for the  $k$ th tessera has the form

$$\hat{J}_k^{\text{screened}} = \frac{\text{erf}(\zeta_k |\vec{r} - \vec{r}_k|)}{|\vec{r} - \vec{r}_k|} \quad (0)$$

and thus

$$L_{\mu\nu}^k = -\int \phi_{\mu}(\vec{r}) \frac{\text{erf}(\zeta_k |\vec{r} - \vec{r}_k|)}{|\vec{r} - \vec{r}_k|} \phi_{\nu}(\vec{r}) d\vec{r} \quad (0)$$

$$V_k^{\text{nuc}} = \sum_{J=1}^{n_A} Z_J \frac{\text{erf}(\zeta_k |\vec{r}_k - \vec{R}_J|)}{|\vec{r}_k - \vec{R}_J|} \quad (0)$$

where Eq. (0) is a generalized version of the nuclear-electron attraction integral, and is sometimes referred to as the Coulomb attenuated nuclear-electron attraction integral. With SWIG/ISWIG, smooth analytical gradients are available for the solvated SCF energy gradient with respect to the nuclear coordinates  $\xi$ . The PCM contribution to the gradient can be obtained by differentiating Eq. (0):

$$\begin{aligned} \Delta G_{\text{solv}}^{(\xi)} &= \frac{1}{2} (\mathbf{V}^{\dagger} \mathbf{q})^{(\xi)} = \frac{1}{2} (\mathbf{R} \mathbf{V}^{\dagger} \mathbf{K}^{-1} \mathbf{V})^{(\xi)} \\ &= \mathbf{q}^{\dagger} \mathbf{V}^{(\xi)} - \frac{1}{2} \mathbf{R}^{-1} \mathbf{q}^{\dagger} \mathbf{K}^{(\xi)} \mathbf{q} \end{aligned} \quad (0)$$

In the derivation we utilize Eqs. (0), (0) and the equation for the derivative of the inverse of a matrix:  $(\mathbf{K}^{-1})^{(\xi)} = -\mathbf{K}^{-1} \mathbf{K}^{(\xi)} \mathbf{K}^{-1}$

### 2.3 Need for GPU Acceleration of PCM

GPUs are characterized as stream processors, and are especially suitable for massively data-parallel computing. To take the most advantage of CPU/GPU heterogeneous programming in quantum chemistry calculations, there are two important considerations: (1) which parts are most time-consuming; (2) which parts are most “stream-friendly”, or benefit most from highly parallelizable architectures. For typical molecules with up to ~15,000 basis functions at the HF/DFT level of theory, building the electron repulsion integrals (ERI) is the most computationally expensive part. However, this cost can be dramatically decreased if we take advantage of the high ratio of compute to data transfer in constructing the ERIs. Therefore, it is a common approach for many developers to accelerate just the ERIs on GPUs and leave the less costly portions of the calculations, e.g. OEs and linear algebra, on CPUs.<sup>55-59</sup>

As mentioned above, the three computationally intensive steps are: 1) constructing  $V^c$ , 2) solving the linear system in Eq. (1), and 3) computing  $\Delta F^s$ . Steps 1) and 3) involve OEIs only, and step 2) is linear algebra. One might think that there was no need to accelerate these PCM components on GPUs, since the cost of the SCF is dominated by the evaluation of two-electron repulsion integrals (ERI). However, in large systems the PCM-related OEIs become comparable in cost to ERIs, and GPU acceleration becomes greatly needed. Figure 3 shows the percentage of time taken by OEI evaluations on CPUs and GPUs for the PCM calculation of a peptide. If only the ERIs are accelerated with GPUs, while the PCM related evaluations remain on CPUs, the PCM portion will take 88% of the total runtime. Therefore, one should expect a considerable speed-up if PCM evaluations are also accelerated with GPUs.

### 3. Methods

Our GPU accelerated PCM and related methods have been implemented within TERACHEM.<sup>57-59</sup> The excited state C-PCM implementation is currently available in a development version of TERACHEM. All other C-PCM functionalities discussed in this article are available in the commercial release v1.9. Timings for this program have been obtained using an Intel Xeon E5-2680 v2 CPU clocked at 2.80GHz with a NVIDIA Tesla P100 GPU. For comparison, ground and excited state PCM calculations have been conducted with the CPU-based electronic structure package Q-CHEM,<sup>14, 52, 53</sup> which uses the same cavity definition and PCM scheme. The restrained ESP (or RESP) charge<sup>60, 61</sup> calculations have been conducted with another CPU based electronic structure package GAMESS-US,<sup>46, 47</sup> since it has the same definition of molecular electrostatic potential (MEP)<sup>33-37</sup> points as TERACHEM. Timings of Q-CHEM and GAMESS-US have been obtained using the faster Intel Xeon E5-2643 CPU clocked at 3.30GHz.

For a performance test of the ground state PCM implementation, we use microsolvated trans-penta-2,4-dieniminium cation ( $C_5H_8N^+$ , trans-PSB3),<sup>62</sup> including up to 186 explicit water molecules (and a dielectric continuum surrounding the microsolvated cluster). Molecular geometries for the microsolvated clusters were obtained using the Amber 12 simulation package<sup>63</sup> by taking snapshots from a classical molecular dynamics simulation. These dynamics simulations used the general AMBER

force field (GAFF)<sup>64, 65</sup> to describe the solute molecule and the three-site transferable intermolecular potential (TIP3P)<sup>66</sup> to describe water molecules. In the PCM calculations, the  $\omega$ PBEh<sup>67</sup> exchange-correlation functional was used with the 6-31G\* basis set.<sup>68</sup> For Fock matrix diagonalization, the Matrix Algebra on GPU and Multicore Architectures (MAGMA)<sup>69-71</sup> library is used for matrices larger than 1500×1500. Smaller Fock matrices are diagonalized with Intel MKL<sup>50</sup> on the CPU.

For performance tests of linear response (LR)<sup>72-74</sup> PCM for TDDFT/TDA, we examine a series of microsolvated umbelliferone molecules with explicit water molecules in a sphere with radius ranging from 2-10Å. The TDDFT/TDA energy and analytical gradient calculations were conducted with the 6-31G\* basis set and the  $\omega$ PBEh<sup>67</sup> functional, which is widely used in TDDFT calculations in order to avoid spurious low-lying charge transfer states.

For performance tests of RESP, the SCF and RESP fitting is conducted on the same set of benchmark molecules for LR-PCM. The HF/6-31G\* level of theory is used, since this choice has been made in the development of many popular force fields<sup>64, 65, 75-77</sup> based on RESP charges. The MEP grid used for fitting the charges is built from 4 layers of Connolly surfaces with radii scaling factor starting from 1.4 and incrementing by 0.2 per layer, with a grid density of 5.0 points/Å<sup>2</sup>.

## 4. Results and Discussion

### 4.1 Acceleration of C-PCM on GPUs

In previous work,<sup>25</sup> we introduced a GPU-accelerated SCF-PCM code where the SCF iterations were accelerated by: (1) computing the PCM related OEIs with a parallel algorithm on GPUs and (2) using a dynamic CG threshold with a RBJ preconditioner. In addition, we have now developed an implementation where the matrix-vector products are built on GPUs without explicit storage of the matrices. This further enhanced the performance for large systems.

Benchmarking results for C-PCM single point energy evaluation are shown in Figure 4. The total run time includes both gas phase DFT evaluations and the PCM terms. The observed scaling of  $O(N^{1.63})$  is very similar to the scaling of the gas phase counterpart. The scaling of the PCM integrals and the linear solver are sub-quadratic and quadratic,

respectively. Based on Eq. (0), the total number of elementary integrals scales as  $O(MN^2)$ , as discussed in section 1. Since  $M$  and  $N$  both grow linearly with the number of atoms for these benchmarking systems (as shown in Figure 4b), the PCM OEI might be expected to scale as  $\sim O(N^3)$ . However, our implementation is observed to scale as  $O(N^{1.58})$  for medium/large sized molecules (with more than 250 atoms and 1500 basis sets). The scaling gains come from exploiting sparsity. The OEI in Eq. (0) are calculated as summations of OEIs over primitive basis functions

$$L_{\mu\nu}^k = \sum_{p=1}^{l_\mu} \sum_{q=1}^{l_\nu} c_{\mu p} c_{\nu q} [\chi_p | \hat{J}_k^{\text{screened}} | \chi_q] \quad (0)$$

where  $c_{\mu p}$  is the contraction coefficient. When we contract  $L_{\mu\nu}^k$  to build the Fock matrix contribution  $\Delta F^S(\mathbf{q})$  in Eq. (0), only significant primitive pairs  $(\chi_p \chi_q)$  are included, as determined by using a Schwartz-like bound<sup>78</sup> with a cutoff,  $\epsilon^{\text{screen}}$ , of  $10^{-12}$  atomic units:

$$[\chi_p | \hat{J}_k^{\text{screened}} | \chi_q] \approx 0 \quad \forall [\chi_p \chi_q | \chi_p \chi_q]^{1/2} < \epsilon^{\text{screen}} \quad (1)$$

When building the ESP,  $V_k^c$  in Eq. (0), the screening criteria is augmented with the corresponding density matrix element:

$$[\chi_p | \hat{J}_k^{\text{screened}} | \chi_q] \approx 0 \quad \forall |P_{\mu\nu}| [\chi_p \chi_q | \chi_p \chi_q]^{1/2} < \epsilon^{\text{screen}} \quad (2)$$

In both cases, only the surviving pair quantities are preloaded to the GPU global memory and participate in the integral calculation.

For small systems, the observed scaling looks even better, but it is mainly because the GPU cores are not saturated with tasks for very small systems. Our linear solver is observed to scale as  $O(N^{1.92})$ , close to the expected  $O(N^2)$  scaling of matrix-vector multiplication in CG. With the acceleration of a single GPU, the C-PCM  $\omega$ PBEh/6-31G\* calculation for the largest system studied in this work (707 atoms, 4495 basis functions) can be completed within 8 minutes.

We previously showed that the GPU-accelerated RBJ preconditioned conjugate gradient method achieved significant speed-ups over both direct inversion and Jacobi preconditioned CG on CPUs. However, as our previous algorithm explicitly stored the  $\mathbf{K}$  matrix in memory, it required  $O(M^2)$  memory which becomes prohibitive for molecules

with  $M > 10^4$  (ca. 200 atoms assuming a Lebedev grid density of 110 pts/atom). To improve upon this, we implemented the matrix-vector product,  $\mathbf{K}\mathbf{q}$ , on GPUs without storing the  $\mathbf{K}$  matrix. This is accomplished in two steps: 1) Pre-computation (on the CPU) of the diagonal elements  $\mathbf{K}_{kk}$  based on Eq. (0) 2) Construction of the off-diagonal elements  $\mathbf{K}_{kl}$  on the GPU simultaneously with evaluation of the matrix-vector product  $\mathbf{K}\mathbf{q}$ , as shown schematically in Figure 5. In detail, each thread is assigned a tessera,  $i$  ( $\mathbf{r}_i$ ,  $\zeta_i$ ), and accumulates contributions to the ESP at this tessera from all other tesserae:  $q_j \text{erf}(\zeta'_{ij} |\mathbf{r}_i - \mathbf{r}_j|) / |\mathbf{r}_i - \mathbf{r}_j|$  for  $i \neq j$ . We pre-load  $(\zeta_j, \mathbf{r}_j)$  to shared memory in a coalesced pattern to reduce memory traffic. Finally, the  $\mathbf{K}\mathbf{q}$  vector is copied back to the CPU and the contributions from diagonal elements of  $\mathbf{K}$  are added in (on the CPU). The whole process requires only  $O(M)$  memory on both GPU and CPU. The standard Jacobi preconditioner is used in this algorithm for simplicity. However, the RBJ and other preconditioners could also be applied with minor modification.

In Figure 6, we compare the GPU-based CG solver with the CPU based RBJ CG solver. For small molecules (less than  $\approx 32$  atoms with grid density of 110 points/atom), the overhead for CPU-GPU data transfer is relatively significant and the CPU based solver is slightly more efficient (Figure 6a). For larger molecules, the GPU solver is significantly faster than the CPU version. Furthermore, we can take advantage of efficient single-precision arithmetic on GPUs, which can significantly accelerate the calculation. For our tests on P100 GPUs, the CG runtime is reduced by about half for large systems. Larger gains will be achieved when using commodity cards that have poor double precision support. Our benchmark data shows that the use of single-precision for the GPU calculation of off-diagonal contributions does not introduce significant errors to the total energy of the solvated system (Figure 6c). In this case, we use double precision for the diagonal elements that are calculated on the CPU. Thus, what we call “single precision” here is actually a mixed-precision scheme. Because of this, the absolute errors in both the total energy and solvation energy introduced by single precision arithmetic are always less than  $10^{-5}$  Hartree even for large molecules with thousands of basis functions and tens of thousands of tesserae. Thus, CG implemented with mixed GPU single precision/CPU double precision will usually be preferred. Due to the significant acceleration achieved in

CG, the ratio of time spent on PCM related evaluations (PCM integrals and CG) relative to total time is greatly reduced (Figure 6d). With the CPU-based CG solver, PCM takes 25%-38% of the total runtime and this percentage increases with system size. The GPU-based mixed-precision solver decreases this ratio to less than 15% for large molecules, significantly outperforming our previously reported<sup>25</sup> typical ratio of 30% for both CPU- and GPU-based implementations.

#### 4.2 Extensions of C-PCM implementation on GPUs: other electronic structure methods

Our C-PCM implementation with GPU acceleration is not only important for ground state HF/DFT, but also lays the foundation for GPU accelerated C-PCM implementation of higher-level quantum chemistry methods. For gas phase quantum chemistry methods, the GPU algorithms for the  $J$  and  $K$  builds (Coulomb and exchange operators, respectively)<sup>56-59</sup> are the basis for the implementation of higher-level methods. GPU acceleration of CIS/TDDFT<sup>79</sup> and CASSCF<sup>80, 81</sup> can be easily achieved through reformulation of the methods in terms of  $J$  and  $K$  builds. Similarly, for C-PCM, the OEI build ( $\Delta\mathbf{F}^S$  and  $V^c$ ) and the linear solver play a pivotal role. Implementation of C-PCM in more complicated quantum chemistry methods can be realized by formulation in terms of these fundamental builds. Such high-level methods might include CIS/TDDFT, CASSCF,<sup>82-86</sup> CASPT2,<sup>87, 88</sup> and newly emergent multi-reference DFT methods like REKS.<sup>89-92</sup> We are especially interested in multi-state versions of the last three methods,<sup>93-97</sup> which are able to treat electronic excited and ground states in a balanced way,<sup>98-101</sup> to correctly describe conical intersections.<sup>102</sup> As an example, here we present our GPU accelerated C-PCM implementation of CIS/TDDFT.

To study the excited state properties of solvated molecules, we implemented both the state-specific (SS)<sup>72, 103</sup> and linear response (LR)<sup>72-74</sup> PCM for TDDFT/TDA, with both equilibrium and non-equilibrium<sup>104, 105</sup> solvation. For SS-PCM, the PCM term only enters the ground state SCF iterations, leaving the working equations of TDDFT virtually unchanged. Only minor modifications need to be performed to separate the “fast” and “slow” polarizations,<sup>104, 105</sup> and an external iteration<sup>103, 106</sup> is performed for TDDFT. The LR-PCM implementation is more complicated, as PCM terms directly enter the TDDFT working equations. These terms also change the coupled-perturbed Kohn-Sham

(CPKS)<sup>107-110</sup> equations needed for analytical gradient evaluation. However, they can be reformulated in terms of the PCM OEI builds for both energy and gradient evaluation. The LR-PCM TDDFT/TDA energies are computed by solving Eq. (16), which represents a solvated version of the gas-phase Casida equation, shown here for TDDFT with the Tamm-Dancoff approximation (TDA):<sup>29-32</sup>

$$\mathbf{A}^{\text{sol}} \mathbf{X} = \omega \mathbf{X} \quad (2)$$

The matrix  $\mathbf{A}^{\text{sol}}$  is the response matrix of the solvated system and defined as the sum of the gas phase response matrix  $\mathbf{A}$  and a PCM correction term  $\mathbf{M}$

$$\mathbf{A}_{ia\sigma, j b\sigma'}^{\text{sol}} = \mathbf{A}_{ia\sigma, j b\sigma'} + \mathbf{M}_{ia\sigma, j b\sigma'} \quad (2)$$

$$\mathbf{M}_{ia\sigma, j b\sigma'} = \sum_{kl} L_{ia\sigma}^k Q(\epsilon_r)_{kl} L_{j b\sigma'}^l \quad (2)$$

Here,  $L_{ia\sigma}^k$  represents the electrostatic potential at tessera  $k$  induced by the charge distribution of the KS orbital pair  $\psi_{i\sigma}(\mathbf{r})\psi_{a\sigma}(\mathbf{r})$ , and is the MO transform of  $L_{\mu\nu}^k$  in Eq. (0). The matrix,  $Q(\epsilon) = \mathbf{K}^{-1}\mathbf{R}$ , is often called the response matrix,<sup>74</sup> because the action of  $Q(\epsilon)$  on  $\mathbf{V}$  produces the surface charges,  $\mathbf{q}$ , that are induced in response to the electrostatic potential  $\mathbf{V}$ . The linear-response dielectric constant is represented by  $\epsilon_r$ , which can represent the solvent static dielectric constant for equilibrium solvation, or the optical dielectric constant for non-equilibrium solvation. In the Davidson<sup>111</sup> iteration of TDDFT/TDA, we only evaluate the matrix-vector product of the response matrix with a trial eigenvector,  $\mathbf{A}^{\text{sol}} \mathbf{X}$ , which leads to a PCM contribution in the AO basis of the form

$$I_{\mu\nu} = \sum_{\gamma\tau} M_{\mu\nu, \gamma\tau} X_{\gamma\tau} = \sum_{\gamma\tau} \sum_{kl} L_{\mu\nu}^k Q(\epsilon_r)_{kl} L_{\gamma\tau}^l X_{\gamma\tau} \quad (2)$$

This term can be reformulated in terms of ground state PCM builds in the following steps:

- (1) Build the ESP induced on the cavity surface by the transition density  $\mathbf{X}$

$$V_k^{e,lr} = \sum_{\gamma\tau} L_{\gamma\tau}^k X_{\gamma\tau} \quad (2)$$

Notice, this step reuses the GPU module for building  $V^e$  of the ground state.

(2) Solve for surface charges in response to  $V_k^{e,lr}$  :

$$\tilde{\mathbf{q}} = \mathbf{Q}(\epsilon_{lr}) \mathbf{V}^{e,lr} \quad (2)$$

i.e.  $\mathbf{K}\tilde{\mathbf{q}} = \mathbf{R}\mathbf{V}^{e,lr}$  with accelerated CG.

(3) Build a 1-electron Fock-like contribution  $I_{\mu\nu} = \sum_k L_{\mu\nu}^k \tilde{q}_k$ . Again, this benefits

from GPU code reuse, exploiting the module for building  $\Delta F_{\mu\nu}^{rS}$ .

To evaluate the analytical gradient of the LR-PCM excited state energy,<sup>112, 113</sup> the PCM correction matrix  $\mathbf{M}$  of Eq. (2), and its partial derivatives with respect to nuclear gradient, enter both the CPKS equation and the analytical gradient of the excited state energy. The detailed formulas are available elsewhere.<sup>112, 113</sup> Here we just describe how these terms can be evaluated with the ground state GPU accelerated PCM module. In the CPKS of LR-PCM, all PCM terms are expressed in terms of the matrix-vector product shown in Eq. (2). For the final gradient evaluation, the derivatives of  $\mathbf{M}$  with respect to nuclear coordinates are also needed:

$$\begin{aligned} & \sum_{\mu\nu,\kappa\lambda} M(\epsilon_{lr})_{\mu\nu,\kappa\lambda}^{(\xi)} X_{\mu\nu} X_{\kappa\lambda} \\ &= \sum_{\mu\nu,\kappa\lambda} \sum_{kl} \left( L_{\mu\nu}^k \right)_{kl}^{(\xi)} Q(\epsilon_{lr})_{kl} L_{\mu\nu}^l X_{\mu\nu} X_{\kappa\lambda} + \sum_{\mu\nu,\kappa\lambda} \sum_{kl} L_{\mu\nu}^k Q(\epsilon_{lr})_{kl} \left( L_{\mu\nu}^l \right)_{kl}^{(\xi)} X_{\mu\nu} X_{\kappa\lambda} \\ & \quad + \sum_{\mu\nu,\kappa\lambda} \sum_{kl} L_{\mu\nu}^k Q(\epsilon_{lr})_{kl}^{(\xi)} L_{\mu\nu}^l X_{\mu\nu} X_{\kappa\lambda} \\ &= 2 \left( \sum_k \tilde{\mathbf{q}}^\dagger \left( \mathbf{V}^{e,lr} \right)^{(\xi)} - \frac{1}{2} \mathbf{R}^{-1} \tilde{\mathbf{q}}^\dagger \mathbf{K}^{(\xi)} \tilde{\mathbf{q}} \right) \end{aligned} \quad (2)$$

The right-hand-side has similar form as the ground state solvation energy gradients in Eq. (0). The only difference is that the ESP,  $\mathbf{V}^{e,lr}$ , and the PCM charge vector,  $\tilde{\mathbf{q}}$ , now correspond to the transition density, as defined in Eq. (2) and (2). Therefore, it can be calculated with the same subroutine used for the ground state solvation energy gradients implementing Eq. (0).

We tested the performance of the GPU based LR-PCM analytical gradient evaluation with a series of microsolvated 7-hydroxycoumarin (Umbelliferone) molecules with explicit water molecules in spheres ranging from 2-10 Å as shown in Figure 7. Umbelliferone ( $C_9H_6O_3$ , Scheme 1) is a well-known fluorophore with high quantum yield<sup>114-116</sup> and has industrial uses as a sunscreen agent.<sup>117</sup> Its lowest electronic excitation ( $S_0 \rightarrow S_1$ ) is a HOMO-LUMO transition consisting of a  $\pi \rightarrow \pi^*$  excitation. The number of basis functions and the number of PCM grid points increase linearly with system size, as shown in Figure 7. Figure 8 shows the total wall time for each calculation of the analytic gradient in microsolvated Umbelliferone. The GPU-based LR-PCM calculation is only negligibly more time-consuming than the GPU-based gas phase TDDFT implementation, and both of these are significantly faster than CPU-based implementations, as demonstrated by comparison with the CPU-based Q-CHEM code. The speed-up compared to the CPU-based Q-CHEM code grows rapidly from about 20x for small molecules to 80x for medium sized molecules. Speed-ups for larger molecules are not shown because the CPU-based calculations are too time-consuming.

#### 4.3 Extensions of C-PCM implementation on GPUs: ESP charge fitting

Apart from the obvious application to solvation models, our GPU-based PCM algorithm also benefits other types of calculations that require evaluation of the ESP. A typical example is ESP-based calculation of atomic charges,<sup>33-37</sup> which are crucial for the development of empirical force fields.<sup>118</sup> The ESP charges are fitted to reproduce the molecular electrostatic potential (MEP)<sup>33-37</sup> at a large number of points defined on three-dimensional surfaces around the molecule of interest. The procedure for building the MEP grid is very similar to building the PCM cavity surface, and evaluating the MEP is equivalent to evaluating  $V^c$  in Eq. (0). Due to the similarities between PCM and ESP evaluation, the GPU based PCM OEI modules can be directly adapted for computing the ESP. The MEP grid discretization schemes<sup>35, 40, 119-121</sup> are often different from the ones used in PCM, and vary for different types of force field development. For the Amber force fields,<sup>65, 122, 123</sup> the MEP is usually calculated on a series of Connolly surfaces.<sup>35, 40</sup> Multiple layers of Connolly surfaces are built outside the vdW surface of the molecule with a constant increment of the radius (Figure 9). Figure 6 shows how the number of

MEP and PCM grid points increase with molecular size for a representative example. The restrained ESP (or RESP) charge model<sup>60, 61</sup> is a variation of ESP charge fitting developed by Kollman and coworkers. It applies weak hyperbolic restraints to address the ill-behaved ESP charges associated with statistically poorly determined centers (“buried” atoms<sup>124, 125</sup>). In the following discussion we focus on the RESP charge fitting implemented in TERACHEM. However, our MEP evaluation algorithm could also be easily adapted to calculate any atomic charges based on the ESP.

We compared the runtime of RESP charge calculation of TERACHEM and GAMESS on the same set of molecules used for benchmarking LR-PCM calculations in the previous section (Figure 10). With the TERACHEM GPU implementation, the total runtime of RESP (including HF/DFT SCF, MEP evaluation and charge fitting) is only 1.2x-1.4x of the time for a pure SCF run. In other words, the extra cost to obtain RESP charges is almost negligible compared to the preceding SCF calculation. This fast RESP implementation is not only useful for force field parameterization, but also applicable to other methods that require calculation of electrostatic interaction energies. Examples include various embedding methods, as well as fragment-based methods<sup>126</sup> such as the recently developed ab initio exciton model.<sup>127</sup>

## 5 Conclusions

We have briefly reviewed the formalism of C-PCM and its implementations in many quantum chemistry packages. We demonstrated the necessity of accelerating PCM on GPUs as a result of the GPU acceleration of the solvent-free calculations. Motivated by this, we developed the techniques to accelerate the SWIG/ISWIG version of C-PCM in our previous work. We further demonstrated our recent implementation of matrix-free CG build on GPUs significantly improves both the memory and the time scaling for PCM calculations of very large molecules. The  $O(N^{1.63})$  time scaling for C-PCM HF/DFT calculations shows that our C-PCM code is one of the most efficient ones available.

Our GPU algorithms are highly transferrable and enable C-PCM implementation for other quantum chemistry methods besides HF and DFT, by exploiting modular PCM OEI and linear equation solver routines. As an example, we discussed the reformulation of PCM-related terms in LR-PCM TDDFT/TDA energy and gradient evaluations. We achieve speedups compared to conventional CPU-based codes of 20-80X for small to

medium-size molecules. We further demonstrated the transferability of our algorithm to methods other than PCM that require ESP evaluations, with the RESP charge derivation as an example. We achieved speedups compared to CPU-based codes of 20-130X for small to medium-size molecules, and the computational overhead introduced by RESP calculation is less than 50%. This enables efficient atomic charge evaluation for force field development and embedding methods.

In the future, we will extend our acceleration strategies to PCM implementation in multi-reference quantum chemistry methods, such as CASSCF, CASPT2, and REKS. Accelerating PCM in these methods will enable excited state *ab initio* molecular dynamics of molecules in solvent environment, thus shedding light on photochemical mechanisms in biological systems.

## **6. Acknowledgement**

The authors acknowledge support by the Department of Energy under grant number DE-SC0018096 (to H.J.K. and F.L). H.J.K. holds a Career Award at the Scientific Interface from the Burroughs Wellcome Fund. TJM is grateful to the Department of Defense (Office of the Assistant Secretary of Defense for Research and Engineering) for a National Security Science and Engineering Faculty Fellowship (NSSEFF). D.M.S. is grateful to the National Science Foundation (NSF) for a graduate fellowship. This work was carried out in part using computational resources from the Extreme Science and Engineering Discovery Environment (XSEDE), which is supported by National Science Foundation grant number ACI-1548562. This work used the XStream computational resource, supported by the National Science Foundation Major Research Instrumentation program (ACI-1429830).

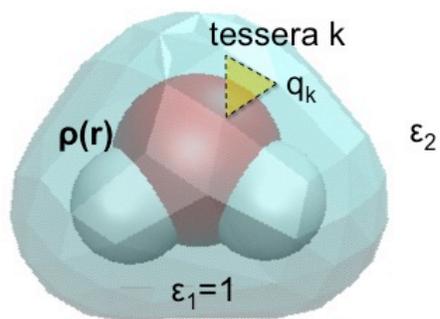
## **Author Information**

Corresponding Author

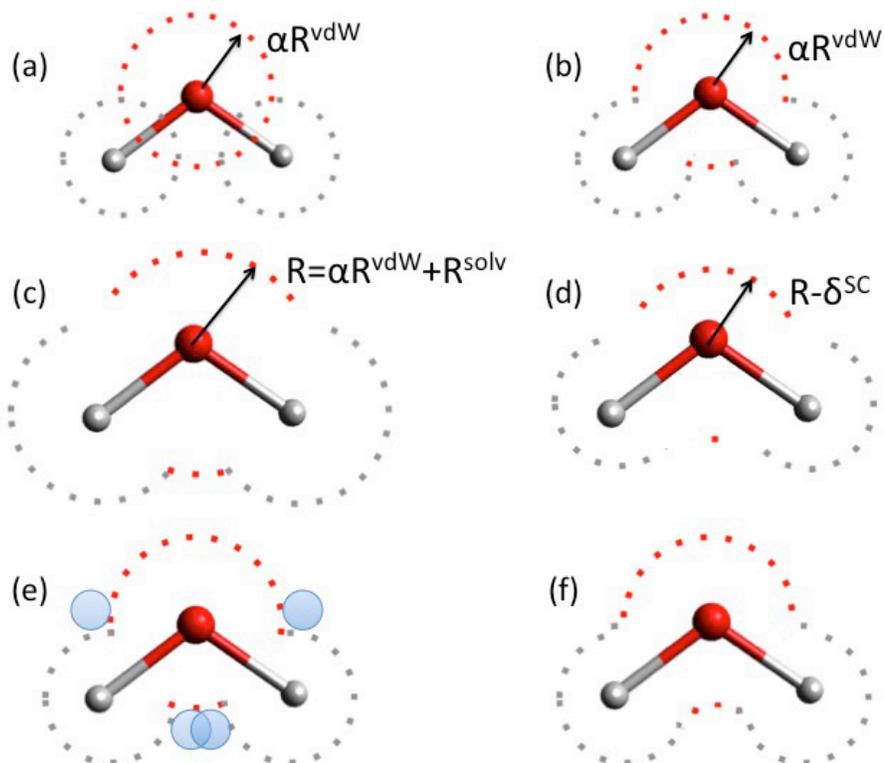
\*E-mail: todd.martinez@stanford.edu

**Table 1.** Combinations of tessellation schemes (rows) and cavity types (columns) in C-PCM implementations for some representative quantum chemistry software packages. Abbreviations used are: GAMESS (GM), Gaussian (GS), NWChem (NW), ORCA (OR), Q-CHEM (QC) and TERACHEM (TC).

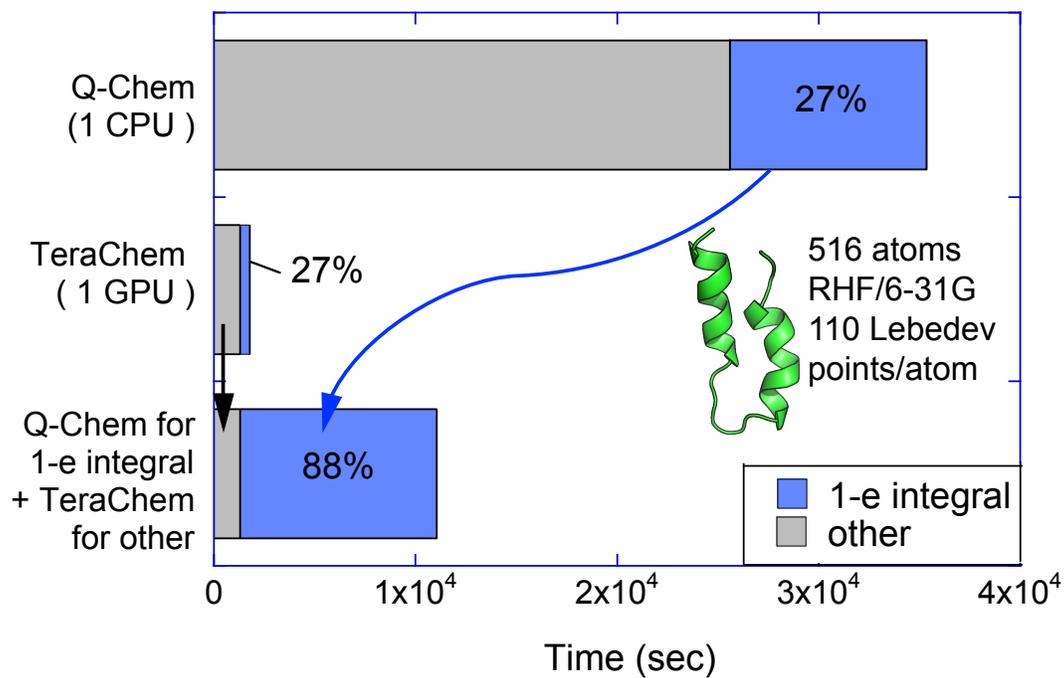
			Scaled VDW/SAS	SES
non-smooth gradient	polyhedron	GEPOL-GB <sup>42-45</sup>	GM,OR	GM,OR
		GEPOL-AS <sup>54</sup>	GM	GM
		GEPOL-RT <sup>54</sup>	GM	GM
		Klamt's <sup>5</sup>	TC	
	Lebedev grid	Fixed	QC, TC	QC
smooth gradient	polyhedron	SWIG <sup>13</sup> (Klamt's <sup>5</sup> )	NW	
		CSC (GEPOL) <sup>15</sup>	GS	GS
	Lebedev grid	SWIG <sup>13</sup>	TC, QC	QC
		ISWIG <sup>14</sup>	TC, QC	QC
	other	FIXPVA <sup>128</sup>	GM	GM



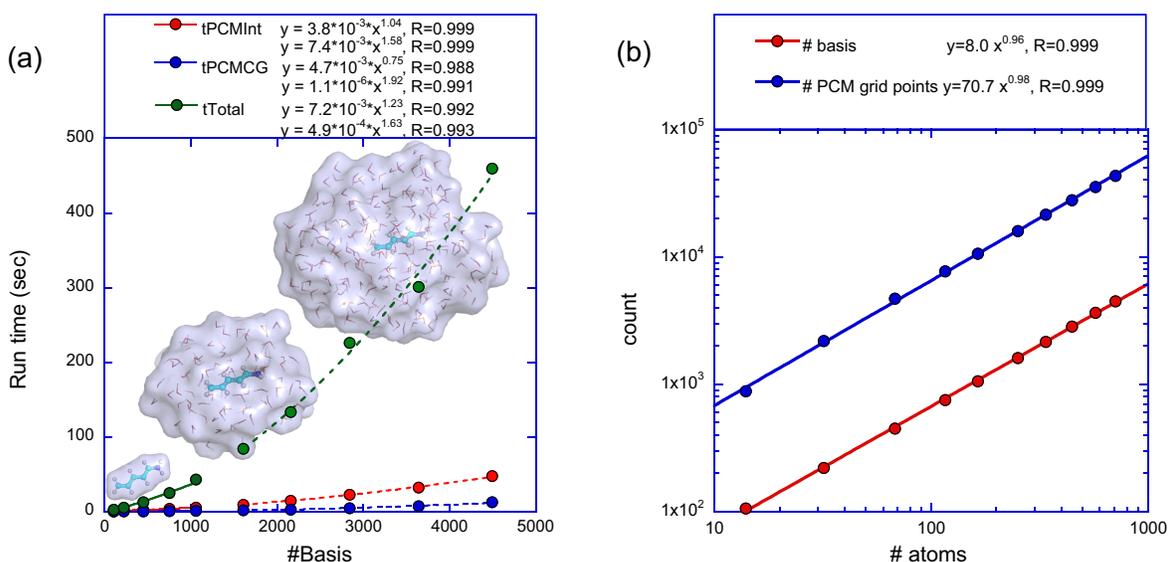
**Figure 1:** Schematic demonstration of polarizable continuum model.



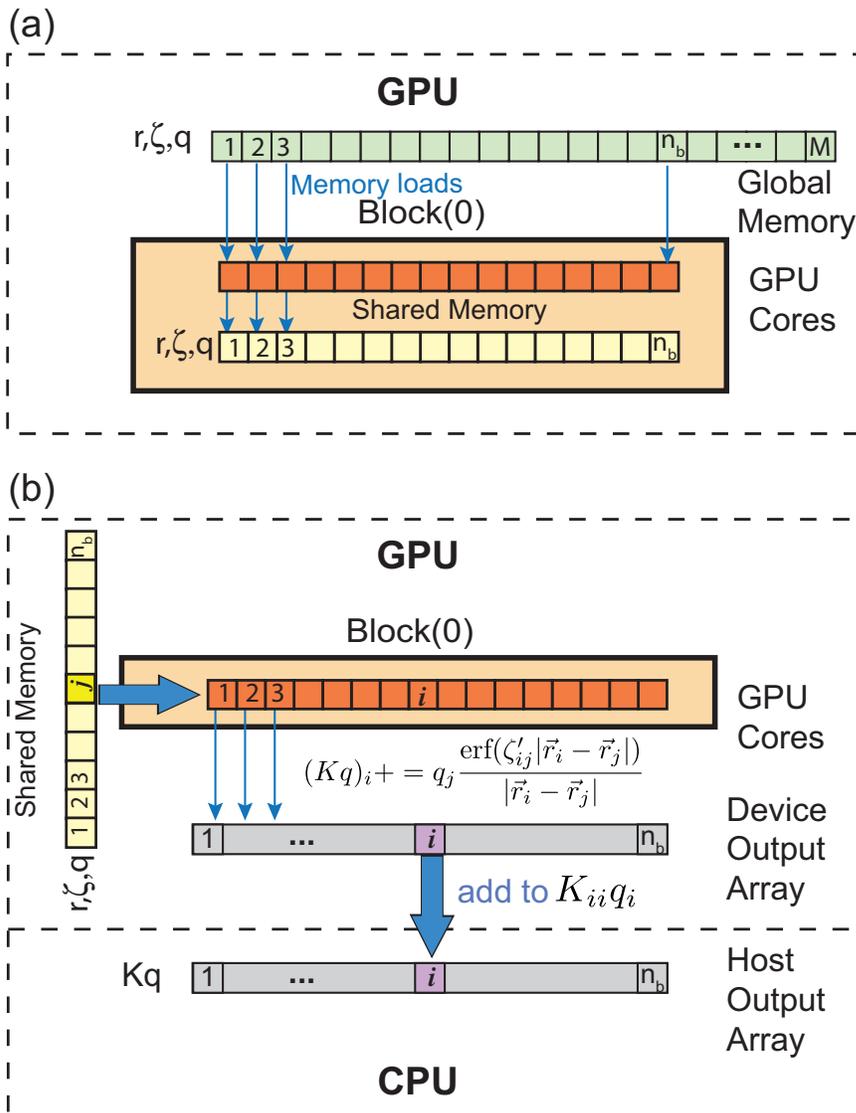
**Figure 2:** The construction of different types of solute cavities (schematic), using a water molecule as an example solute. (a) Interlocking spheres with scaled van der Waals radii centered on the solute atoms (scaling factor  $\alpha \geq 1$ ). (b) Remove the surface buried inside the cavity and form the scaled van der Waals surface. This is the currently used definition of solvent accessible surface (SAS). (c) Form the surface built in (b), expand by the solvent radii  $R^{solv}$ . (d) Contract the radii of surface built in (c) by  $\delta^{SC}$  and form the Klamt solvent accessible surface (SAS(K)) in the original COSMO method. (e) Roll a probe (representing the solvent) around the SAS in (b). (f) Take the area explored by center of the probe as the solvent excluded surface (SES, or Connolly surface).



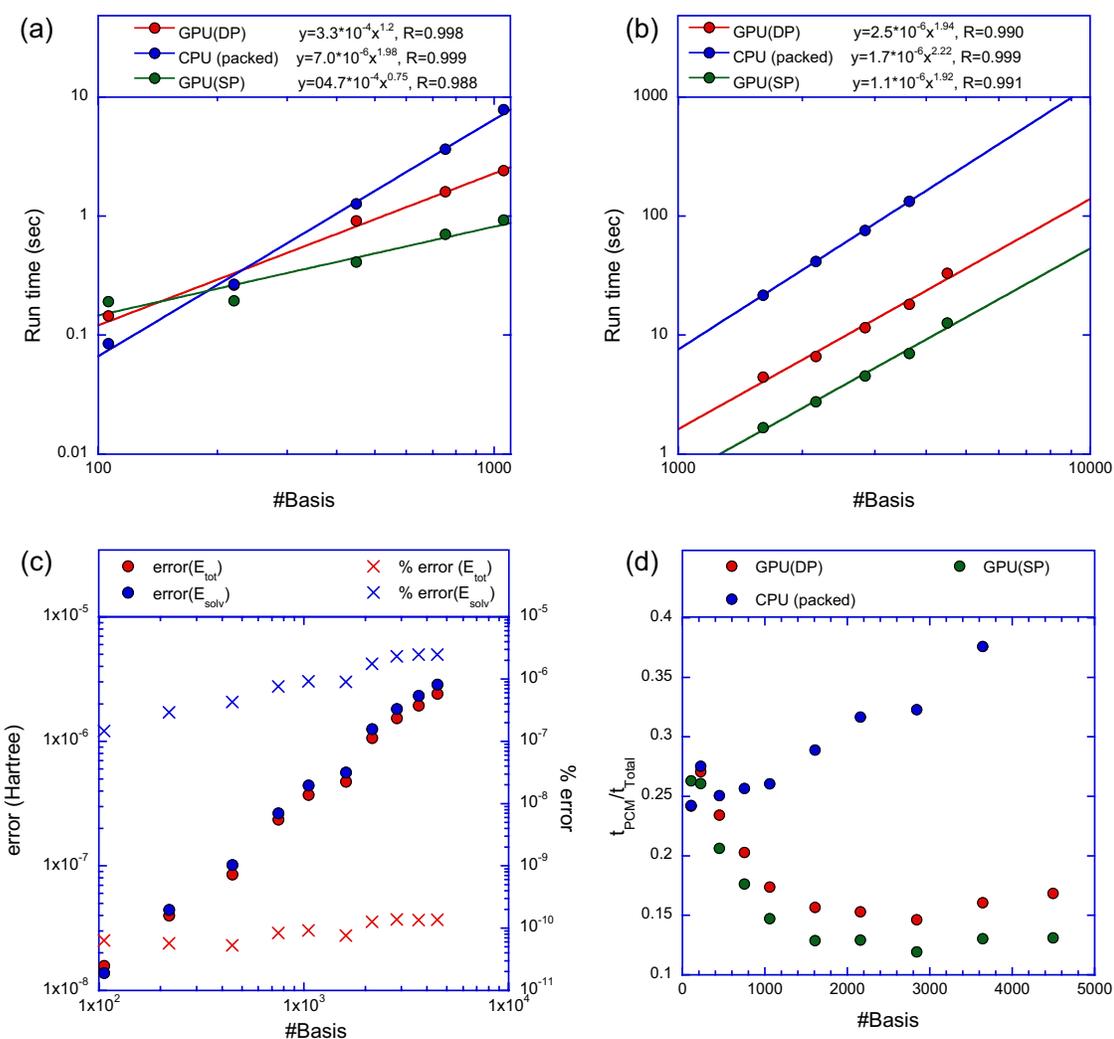
**Figure 3:** Percentage of time taken by OEIs in a polarizable continuum model calculation of a peptide (PDB ID: 2KJM) with **TERA-CHEM** and **Q-CHEM**.



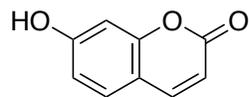
**Figure 4:** Timing for C-PCM  $\omega$ PBEh/6-31G\* single point energy evaluation for trans-PSB3 and a series of microsolvated trans-PSB3 molecules. (a) The total timing involves the whole SCF calculation (both gas phase portion and PCM related). The time for PCM integrals includes calculating the electrostatic potential and building the Fock matrix contribution. The structures for trans-PSB3 molecule and its microsolvated forms in water spheres of radii 6 Å and 10 Å are shown with respective cavity surfaces in the inset. C, N, O, H atoms are colored with cyan, blue, red, and white, respectively. Timings were obtained using 1 NVIDIA Tesla P100 GPU accelerator. For the run time, a power law fit is conducted separately for small-medium systems (up to 1067 basis functions) and larger systems. Intel MKL and GPU-based MAGMA libraries are used for the Fock matrix diagonalization, respectively. (b) Number of basis functions and number of PCM grid points as a function of the number of atoms in the benchmark molecules.



**Figure 5:** Algorithm for building  $\mathbf{Kq}$  product on GPUs. At the top of the graph, the pale green array represents data for each tessera (Gaussian charge distribution), including the center coordinates, Gaussian exponents, and charge. The GPU cores are represented by orange squares (threads) embedded in pale yellow rectangles (1 dimensional blocks with  $n_b=16$  threads/block). The algorithm includes two steps. (a) Each thread fetches the data of the tessera corresponding to its thread index and copies it to shared memory as a coalesced transfer. (b) Each thread is assigned to the data of a tessera corresponding to its ID,  $i$ , and loops over all tesserae  $j$  loaded in the shared memory to calculate the contribution to the  $i$ -th entry of the  $\mathbf{Kq}$  product:  $K_{ij}q_j$ . Since the number of tesserae is usually larger than the block size, the tessera data will be loaded in batches into shared memory, and the two steps will be repeated until  $j$  loops over all tesserae. We show only the 0th block. For block  $b$ , the tessera indices  $i$  corresponding to each thread in the block are  $b*n_b, b*n_b+1, \dots, (b+1)*n_b-1$

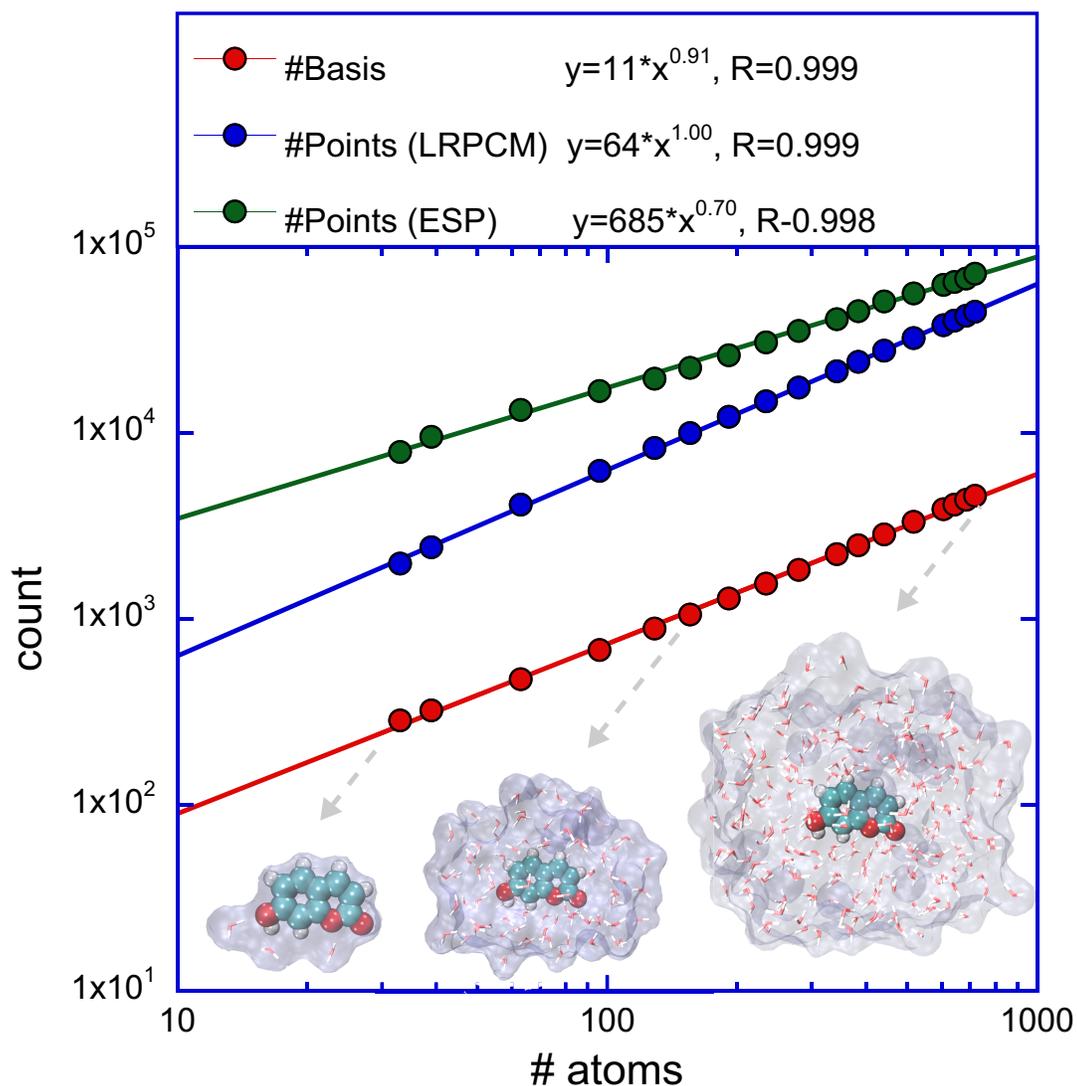


**Figure 6:** Timing for different TERACHEM PCM CG solver implementations for trans-PSB3 and a series of microsolvated trans-PSB3 molecules. The CPU (packed) RBJ solver explicitly stores the upper-triangle of the PCM matrix and uses MKL on the CPU to perform the matrix-vector product. GPU(DP) and GPU(SP) refer to forming the matrix-vector product on the GPU (without explicit storage of the matrix) using double precision (DP) and single precision (SP). Empirical scaling of run time is obtained from a power law fit. The fit is done separately for (a) small/medium (up to 1067 basis functions) and (b) large systems, respectively. (c) Absolute error in total energy of the solvated system ( $E_{\text{tot}}$ ) and solvation energy ( $E_{\text{solv}}$ ) calculated with GPU(SP) compared with GPU(DP). (d) Ratio of time spent on PCM related evaluations (PCM integrals and CG solver) with respect to the total run time (PCM + gas phase).

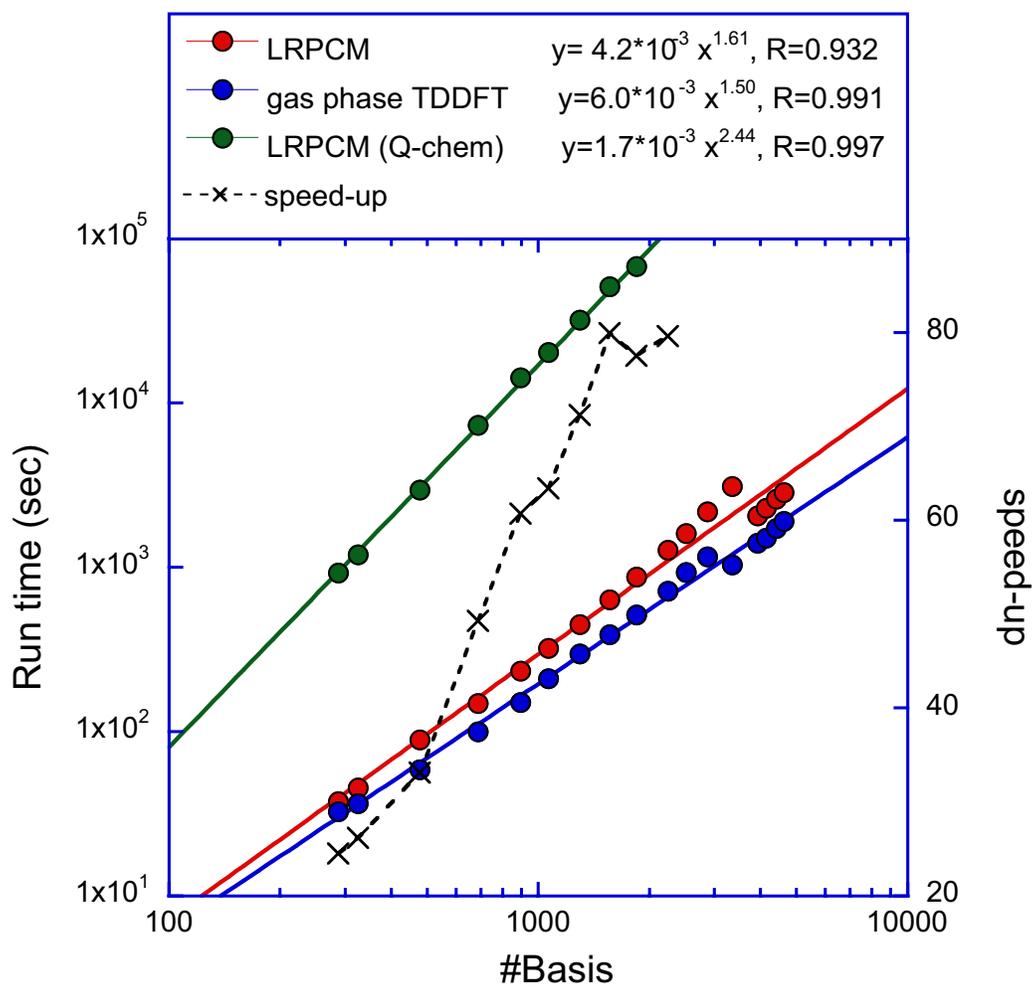


Umbelliferone

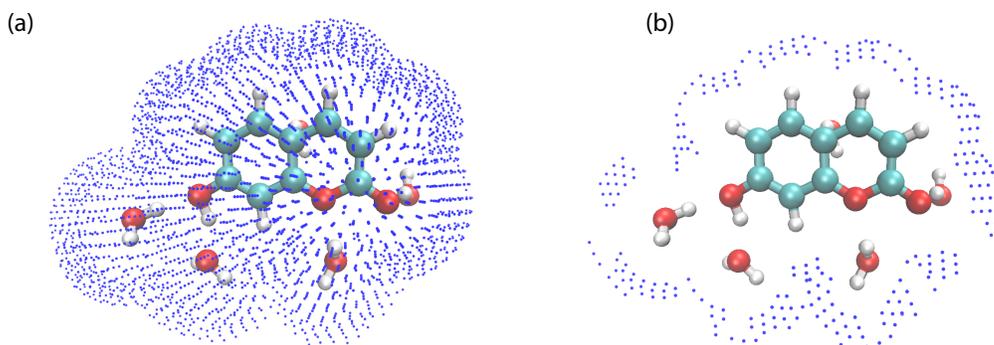
**Scheme 1:** Chemical structure of Umbelliferone (IUPAC name: 7-Hydroxychromen-2-one)



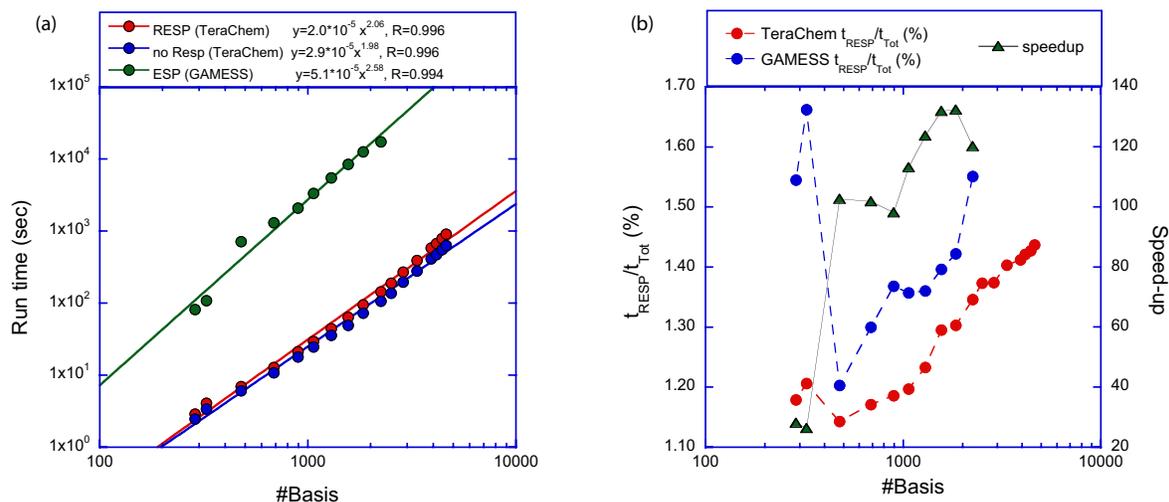
**Figure 7:** Systems used for benchmarking LR-PCM and RESP calculations: a series of microsolvated umbelliferone molecules with increasing numbers of explicit water molecules. Red dots show the number of basis functions in each system (6-31G\* basis set). Blue dots show the number of cavity surface points used for LRPCM benchmarking, with Lebedev grid density of 110 points/atom and cavity radii scaling factor of 1.2. Green dots show the number of ESP grid points used for RESP charge benchmarking, with 4 layers of Connolly surfaces with radii scaling factor starting from 1.4 and an increment of 0.2 per layer. Grid density is 5.0 points/Å<sup>2</sup>.



**Figure 8:** Timings for TDDFT/TDA  $\omega$ PBEh/6-31G\* gradient evaluation in gas phase and in LR-PCM equilibrium solvation for a series of microsolvated umbelliferone molecules. Red and blue dots show timings for TERACHEM using 1 NVIDIA Tesla P100 GPU accelerator. Green dots are the timings for Q-CHEM LR-PCM calculations using 1 Intel Xeon E5-2643 CPU clocked at 3.30GHz. The speed-up for TERACHEM vs Q-CHEM for LRPCM is shown by black Xs.



**Figure 9** (a) RESP grid points for microhydrated umbelliferone. MEP grid points were built from 4 layers of Connolly surface with radii scaling factor starting from 1.4 and an increment of 0.2 per layer. Grid density is  $5.0 \text{ points}/\text{\AA}^2$ . The molecules are represented with ball and sticks, with C, O, H colored as cyan, red, and white, respectively. Blue dots represent grid points. (b) Slice of the grid points in the plane of the aromatic ring, highlighting the four-layer structure of the RESP grid.



**Figure 10:** Timings for HF/6-31G\* RESP charge calculation for a series of microhydrated umbelliferone molecules. (a) Red and blue dots represent the timings for TERAChem using 1 NVIDIA Tesla P100 GPU accelerator. Green dots are the timings for GAMESS LR-PCM calculations using 1 Intel Xeon E5-2643 CPU clocked at 3.30GHz. RESP timing includes both SCF and the following RESP charge evaluation. (b) The percentage of time taken by RESP charge evaluation in an SCF calculation followed by RESP charge evaluation is shown with red dots and blue dots, for TERAChem and GAMESS, respectively. The speed-up of TERAChem compared to GAMESS for the RESP charge evaluation is shown as blue dots.

## References

- <sup>1</sup> I. S. Ufimtsev, N. Luehr, and T. J. Martinez, *J. Phys. Chem. Lett.* **2** (2011) 1789.
- <sup>2</sup> B. Mennucci, *WIREs: Comput. Mol. Sci.* **2** (2012) 386.
- <sup>3</sup> J. Tomasi, and M. Persico, *Chem. Rev.* **94** (1994) 2027.
- <sup>4</sup> S. Miertuš, E. Scrocco, and J. Tomasi, *Chem. Phys.* **55** (1981) 117.
- <sup>5</sup> A. Klamt, and G. Schuurmann, *J. Chem. Soc. Perkin Trans. 2* (1993) 799.
- <sup>6</sup> J. Andzelm, C. Kölmel, and A. Klamt, *J. Chem. Phys.* **103** (1995) 9312.
- <sup>7</sup> A. Klamt, *WIREs: Comput. Mol. Sci.* **8** (2017) e1338.
- <sup>8</sup> V. Barone, and M. Cossi, *J. Phys. Chem. A* **102** (1998) 1995.
- <sup>9</sup> T. N. Truong, and E. V. Stefanovich, *Chem. Phys. Lett.* **240** (1995) 253.
- <sup>10</sup> E. Cancès, B. Mennucci, and J. Tomasi, *J. Chem. Phys.* **107** (1997) 3032.
- <sup>11</sup> B. Mennucci, E. Cancès, and J. Tomasi, *J. Phys. Chem. B* **101** (1997) 10506.
- <sup>12</sup> J. Tomasi, B. Mennucci, and E. Cancès, *J. Mol. Struct: THEOCHEM* **464** (1999) 211.
- <sup>13</sup> D. M. York, and M. Karplus, *J. Phys. Chem. A* **103** (1999) 11060.
- <sup>14</sup> A. W. Lange, and J. M. Herbert, *J. Chem. Phys.* **133** (2010) 244111.
- <sup>15</sup> G. Scalmani, and M. J. Frisch, *J. Chem. Phys.* **132** (2010) 114110.
- <sup>16</sup> M. Orozco, and F. J. Luque, *Chem. Rev.* **100** (2000) 4187.
- <sup>17</sup> L. Greengard, and V. Rokhlin, *J. Comp. Phys.* **73** (1987) 325.
- <sup>18</sup> P. Li, H. Johnston, and R. Krasny, *J. Comp. Phys.* **228** (2009) 3858.
- <sup>19</sup> G. Scalmani *et al.*, *Theo. Chem. Acc.* **111** (2004) 90.
- <sup>20</sup> B. Stamm *et al.*, *J. Chem. Phys.* **144** (2016) 054101.
- <sup>21</sup> F. Lipparini *et al.*, *J. Phys. Chem. Lett.* **5** (2014) 953.
- <sup>22</sup> E. Cancès, Y. Maday, and B. Stamm, *J. Chem. Phys.* **139** (2013) 054111.
- <sup>23</sup> F. Lipparini *et al.*, *J. Chem. Theo. Comp.* **9** (2013) 3637.
- <sup>24</sup> F. Lipparini *et al.*, *J. Chem. Phys.* **141** (2014) 184108.
- <sup>25</sup> F. Liu *et al.*, *J. Chem. Theo. Comp.* **11** (2015) 3131.
- <sup>26</sup> G. H. Golub, and C. F. Van Loan, *Matrix computations* (Johns Hopkins University Press, Baltimore, 2013), 4th edn., Johns Hopkins studies in the mathematical sciences.,

- <sup>27</sup> A. Szabo, and N. S. Ostlund, *Modern Quantum Chemistry* (Dover, New York, 1996),
- <sup>28</sup> W. Kohn, and L. J. Sham, *Phys. Rev.* **140** (1965) A1133.
- <sup>29</sup> S. Hirata, and M. Head-Gordon, *Chem. Phys. Lett.* **302** (1999) 375.
- <sup>30</sup> S. Hirata, M. Head-Gordon, and R. J. Bartlett, *J. Chem. Phys.* **111** (1999) 10774.
- <sup>31</sup> A. Dreuw, and M. Head-Gordon, *Chem. Rev.* **105** (2005) 4009.
- <sup>32</sup> F. Furche, and R. Ahlrichs, *J. Chem. Phys.* **117** (2002) 7433.
- <sup>33</sup> F. A. Momany, *J. Phys. Chem.* **82** (1978) 592.
- <sup>34</sup> S. Cox, and D. Williams, *J. Comp. Chem.* **2** (1981) 304.
- <sup>35</sup> U. C. Singh, and P. A. Kollman, *J. Comp. Chem.* **5** (1984) 129.
- <sup>36</sup> L. E. Chirlian, and M. M. Francl, *J. Comp. Chem.* **8** (1987) 894.
- <sup>37</sup> D. E. Williams, *Rev. Comp. Chem.* **2** (2007) 219.
- <sup>38</sup> L. E. McMurchie, and E. R. Davidson, *J. Comp. Phys.* **26** (1978) 218.
- <sup>39</sup> D. M. Chipman, *J. Chem. Phys.* **112** (2000) 5558.
- <sup>40</sup> M. L. Connolly, *J. Appl. Cryst.* **16** (1983) 548.
- <sup>41</sup> V. I. Lebedev, *USSR Comp. Math. Math. Phys.* **16** (1976) 10.
- <sup>42</sup> J. L. Pascual - Ahuir, E. Silla, and I. Tunon, *J. Comp. Chem.* **15** (1994) 1127.
- <sup>43</sup> J. L. Pascual - Ahuir, and E. Silla, *J. Comp. Chem.* **11** (1990) 1047.
- <sup>44</sup> E. Silla, I. Tunon, and J. L. Pascual - Ahuir, *J. Comp. Chem.* **12** (1991) 1077.
- <sup>45</sup> E. Silla *et al.*, *J. Mol. Graph.* **8** (1990) 168.
- <sup>46</sup> M. W. Schmidt *et al.*, *J. Comp. Chem.* **14** (1993) 1347.
- <sup>47</sup> M. S. Gordon, and M. W. Schmidt, in *Theory and applications of computational chemistry* (Elsevier, 2005), pp. 1167.
- <sup>48</sup> M. J. Frisch *et al.*, Wallingford, CT, 2016).
- <sup>49</sup> F. Neese, *WIREs: Comput. Mol. Sci.* **2** (2012) 73.
- <sup>50</sup> O. Treutler, and R. Ahlrichs, *J. Chem. Phys.* **102** (1995) 346.
- <sup>51</sup> C. W. Murray, N. C. Handy, and G. J. Laming, *Mol. Phys.* **78** (1993) 997.
- <sup>52</sup> Y. Shao *et al.*, *Mol. Phys.* **113** (2015) 184.
- <sup>53</sup> A. W. Lange, and J. M. Herbert, *J. Chem. Phys.* **134** (2011) 117102.
- <sup>54</sup> H. Li, and J. H. Jensen, *J. Comp. Chem.* **25** (2004) 1449.
- <sup>55</sup> K. Yasuda, *J. Comp. Chem.* **29** (2008) 334.
- <sup>56</sup> I. S. Ufimtsev, and T. J. Martinez, *Comp Sci. Eng.* **10** (2008) 26.

- <sup>57</sup> I. S. Ufimtsev, and T. J. Martinez, *J. Chem. Theo. Comp.* **4** (2008) 222.
- <sup>58</sup> I. S. Ufimtsev, and T. J. Martinez, *J. Chem. Theo. Comp.* **5** (2009) 1004.
- <sup>59</sup> I. S. Ufimtsev, and T. J. Martinez, *J. Chem. Theo. Comp.* **5** (2009) 2619.
- <sup>60</sup> C. I. Bayly *et al.*, *J. Phys. Chem.* **97** (1993) 10269.
- <sup>61</sup> W. D. Cornell *et al.*, *J. Am. Chem. Soc.* **115** (1993) 9620.
- <sup>62</sup> M. Garavelli *et al.*, *J. Am. Chem. Soc.* **119** (1997) 6891.
- <sup>63</sup> D. A. Case *et al.*, (University of California, 2012).
- <sup>64</sup> J. Wang *et al.*, *J. Mol. Graph. Model.* **25** (2006) 247.
- <sup>65</sup> J. Wang *et al.*, *J. Comp. Chem.* **25** (2004) 1157.
- <sup>66</sup> W. L. Jorgensen *et al.*, *J. Chem. Phys.* **79** (1983) 926.
- <sup>67</sup> M. A. Rohrdanz, K. M. Martins, and J. M. Herbert, *J. Chem. Phys.* **130** (2009)
- <sup>68</sup> Ditchfie.R, W. J. Hehre, and J. A. Pople, *J. Chem. Phys.* **54** (1971) 724.
- <sup>69</sup> S. Tomov, J. Dongarra, and M. Baboulin, *Par. Comp.* **36** (2010) 232.
- <sup>70</sup> S. Tomov *et al.*, in *Parallel & Distributed Processing, Workshops and Phd Forum (IPDPSW), 2010 IEEE International Symposium on* (IEEE, 2010), pp. 1.
- <sup>71</sup> J. Dongarra *et al.*, in *Numerical Computations with GPUs* (Springer, 2014), pp. 3.
- <sup>72</sup> R. Cammi *et al.*, *J. Chem. Phys.* **122** (2005) 104513.
- <sup>73</sup> G. Scalmani *et al.*, *J. Chem. Phys.* **124** (2006) 094107.
- <sup>74</sup> M. Cossi, and V. Barone, *J. Chem. Phys.* **115** (2001) 4708.
- <sup>75</sup> W. D. Cornell *et al.*, *J. Am. Chem. Soc.* **117** (1995) 5179.
- <sup>76</sup> R. J. Woods *et al.*, *J. Phys. Chem.* **99** (1995) 3832.
- <sup>77</sup> K. N. Kirschner *et al.*, *J. Comp. Chem.* **29** (2008) 622.
- <sup>78</sup> J. L. Whitten, *J. Chem. Phys.* **58** (1973) 4496.
- <sup>79</sup> C. M. Isborn *et al.*, *J. Chem. Theo. Comp.* **7** (2011) 1814.
- <sup>80</sup> E. G. Hohenstein *et al.*, *J. Chem. Phys.* **142** (2015)
- <sup>81</sup> J. W. Snyder *et al.*, *J. Chem. Phys.* **143** (2015)
- <sup>82</sup> H. B. Schlegel, and M. A. Robb, *Chem. Phys. Lett.* **93** (1982) 43.
- <sup>83</sup> D. Hegarty, and M. A. Robb, *Mol. Phys.* **38** (1979) 1795.
- <sup>84</sup> N. Yamamoto *et al.*, *Chem. Phys. Lett.* **250** (1996) 373.
- <sup>85</sup> R. H. Eade, and M. A. Robb, *Chem. Phys. Lett.* **83** (1981) 362.
- <sup>86</sup> M. Frisch *et al.*, *Chem. Phys. Lett.* **189** (1992) 524.

- <sup>87</sup> K. Andersson *et al.*, *J. Phys. Chem.* **94** (1990) 5483.
- <sup>88</sup> K. Andersson, P. Å. Malmqvist, and B. O. Roos, *J. Chem. Phys.* **96** (1992) 1218.
- <sup>89</sup> M. Filatov, and S. Shaik, *Chem. Phys. Lett.* **304** (1999) 429.
- <sup>90</sup> I. d. P. R. Moreira *et al.*, *J. Chem. Theo. Comp.* **3** (2007) 764.
- <sup>91</sup> M. Filatov, *WIREs: Comput. Mol. Sci.* **5** (2015) 146.
- <sup>92</sup> M. Filatov, in *Density-functional methods for excited states*, edited by N. Ferre, M. Filatov, and M. Huix-Rotllant (Springer, Heidelberg, 2016), pp. 97.
- <sup>93</sup> J. W. Snyder Jr, R. M. Parrish, and T. J. Martínez, *J. Phys. Chem. Lett.* **8** (2017) 2432.
- <sup>94</sup> R. N. Diffenderfer, and D. R. Yarkony, *J. Phys. Chem.* **86** (1982) 5098.
- <sup>95</sup> J. Finley *et al.*, *Chem. Phys. Lett.* **288** (1998) 299.
- <sup>96</sup> T. Shiozaki *et al.*, *J. Chem. Phys.* **135** (2011) 081106.
- <sup>97</sup> M. Filatov, F. Liu, and T. J. Martínez, *J. Chem. Phys.* **147** (2017) 034113.
- <sup>98</sup> B. G. Levine *et al.*, *Mol. Phys.* **104** (2006) 1039.
- <sup>99</sup> M. Huix-Rotllant *et al.*, *J. Chem. Theo. Comp.* **9** (2013) 3917.
- <sup>100</sup> S. L. Li *et al.*, *J. Phys. Chem. Lett.* **5** (2014) 322.
- <sup>101</sup> S. Gozem *et al.*, *J. Chem. Theo. Comp.* **8** (2012) 4069.
- <sup>102</sup> D. R. Yarkony, *Chem. Rev.* **112** (2011) 481.
- <sup>103</sup> R. Improta *et al.*, *J. Chem. Phys.* **125** (2006) 054103.
- <sup>104</sup> M. Cossi, and V. Barone, *J. Chem. Phys.* **112** (2000) 2427.
- <sup>105</sup> M. Cossi, and V. Barone, *J. Phys. Chem. A* **104** (2000) 10614.
- <sup>106</sup> R. Improta *et al.*, *J. Chem. Phys.* **127** (2007) 074504.
- <sup>107</sup> R. E. Stratmann, G. E. Scuseria, and M. J. Frisch, *J. Chem. Phys.* **109** (1998) 8218.
- <sup>108</sup> H. Weiss, R. Ahlrichs, and M. Häser, *J. Chem. Phys.* **99** (1993) 1262.
- <sup>109</sup> A. Görling *et al.*, *J. Chem. Phys.* **110** (1999) 2785.
- <sup>110</sup> S. Van Gisbergen, J. Snijders, and E. Baerends, *Comp. Phys. Comm.* **118** (1999) 119.
- <sup>111</sup> E. R. Davidson, *J. Comp. Phys.* **17** (1975) 87.
- <sup>112</sup> J. Liu, and W. Liang, *J. Chem. Phys.* **138** (2013) 024101.
- <sup>113</sup> Y. Wang, and H. Li, *J. Chem. Phys.* **133** (2010) 034108.
- <sup>114</sup> J.-J. Aaron *et al.*, *J. Fluorescence* **5** (1995) 337.
- <sup>115</sup> T. Moriya, *Bull. Chem. Soc. Japan* **61** (1988) 1873.
- <sup>116</sup> K. Azuma *et al.*, *Photochem. Photobio. Sci.* **2** (2003) 443.

- <sup>117</sup> B. Ramesh, and K. Pugalendi, *Life sciences* **79** (2006) 306.
- <sup>118</sup> F.-Y. Dupradeau *et al.*, *Phys. Chem. Chem. Phys.* **12** (2010) 7821.
- <sup>119</sup> C. M. Breneman, and K. B. Wiberg, *J. Comp. Chem.* **11** (1990) 361.
- <sup>120</sup> M. A. Spackman, *J. Comp. Chem.* **17** (1996) 1.
- <sup>121</sup> R. H. Henchman, and J. W. Essex, *J. Comp. Chem.* **20** (1999) 483.
- <sup>122</sup> J. A. Maier *et al.*, *J. Chem. Theo. Comp.* **11** (2015) 3696.
- <sup>123</sup> V. Hornak *et al.*, *Proteins: Structure, Function, and Bioinformatics* **65** (2006) 712.
- <sup>124</sup> D. E. Williams, *Biopolymers* **29** (1990) 1367.
- <sup>125</sup> K. M. Merz, *J. Comp. Chem.* **13** (1992) 749.
- <sup>126</sup> A. Goez, and J. Neugebauer, in *Frontiers of Quantum Chemistry* (Springer, 2018), pp. 139.
- <sup>127</sup> X. Li *et al.*, *J. Chem. Theo. Comp.* **13** (2017) 3493.
- <sup>128</sup> P. Su, and H. Li, *J. Chem. Phys.* **130** (2009) 074109.