

Linear scaling computation of the Fock matrix. VI. Data parallel computation of the exchange-correlation matrix

Chee Kwan Gan and Matt Challacombe

Citation: *The Journal of Chemical Physics* **118**, 9128 (2003); doi: 10.1063/1.1568734

View online: <http://dx.doi.org/10.1063/1.1568734>

View Table of Contents: <http://scitation.aip.org/content/aip/journal/jcp/118/20?ver=pdfcov>

Published by the [AIP Publishing](#)

Articles you may be interested in

[Linear-scaling implementation of molecular response theory in self-consistent field electronic-structure theory](#)
J. Chem. Phys. **126**, 154108 (2007); 10.1063/1.2715568

[Parallel algorithm for the computation of the Hartree-Fock exchange matrix: Gas phase and periodic parallel ONX](#)
J. Chem. Phys. **125**, 104110 (2006); 10.1063/1.2222359

[Linear scaling computation of the Fock matrix. VII. Parallel computation of the Coulomb matrix](#)
J. Chem. Phys. **121**, 6608 (2004); 10.1063/1.1790891

[Linear scaling computation of the Fock matrix. IV. Multipole accelerated formation of the exchange matrix](#)
J. Chem. Phys. **111**, 6223 (1999); 10.1063/1.479926

[Linear scaling computation of the Fock matrix](#)
J. Chem. Phys. **106**, 5526 (1997); 10.1063/1.473575



Linear scaling computation of the Fock matrix. VI. Data parallel computation of the exchange-correlation matrix

Chee Kwan Gan and Matt Challacombe^{a)}

Theoretical Division, Group T-12, MS B268, Los Alamos National Laboratory, Los Alamos, New Mexico 87545

(Received 4 November 2002; accepted 27 February 2003)

Recently, early onset linear scaling computation of the exchange-correlation matrix has been achieved using hierarchical cubature [J. Chem. Phys. **113**, 10037 (2000)]. Hierarchical cubature differs from other methods in that the integration grid is adaptive and purely Cartesian, which allows for a straightforward domain decomposition in parallel computations; the volume enclosing the entire grid may be simply divided into a number of nonoverlapping boxes. In our data parallel approach, each box requires only a fraction of the total density to perform the necessary numerical integrations due to the finite extent of Gaussian-orbital basis sets. This inherent data locality may be exploited to reduce communications between processors as well as to avoid memory and copy overheads associated with data replication. Although the hierarchical cubature grid is Cartesian, naive boxing leads to irregular work loads due to strong spatial variations of the grid and the electron density. In this paper we describe equal time partitioning, which employs time measurement of the smallest sub-volumes (corresponding to the primitive cubature rule) to load balance grid-work for the next self-consistent-field iteration. After start-up from a heuristic center of mass partitioning, equal time partitioning exploits smooth variation of the density and grid between iterations to achieve load balance. With the 3-21G basis set and a medium quality grid, equal time partitioning applied to taxol (62 heavy atoms) attained a speedup of 61 out of 64 processors, while for a 110 molecule water cluster at standard density it achieved a speedup of 113 out of 128. The efficiency of equal time partitioning applied to hierarchical cubature improves as the grid work per processor increases. With a fine grid and the 6-311G(df,p) basis set, calculations on the 26 atom molecule α -pinene achieved a parallel efficiency better than 99% with 64 processors. For more coarse grained calculations, superlinear speedups are found to result from reduced computational complexity associated with data parallelism. © 2003 American Institute of Physics.

[DOI: 10.1063/1.1568734]

I. INTRODUCTION

Density functional theory (DFT) and its variant, the hybrid Hartree–Fock/density functional theory (HF/DFT) have proven to be accurate and computationally efficient. These model chemistries are routinely used by conventional Gaussian-orbital quantum chemistry codes to excellent advantage. Recently, significant progress has been achieved in the algorithmic development of $O(N)$ alternatives to these conventional methods, where compute times scale linearly with system size N .^{1,2} These linear scaling methods overcome a number of bottlenecks of order $O(N^{2-3})$ associated with conventional methods, including computation of the exact Hartree–Fock exchange matrix,^{3–8} the Coulomb matrix,^{9–14} the exchange-correlation matrix,^{15–18} and density matrix alternatives to an eigensolution of the self-consistent-field (SCF) equations.^{19–30}

The existence of linear scaling methods has been argued from the concept of “nearsightedness” or quantum locality,^{31,32} which implies that the density matrix $n(\mathbf{r},\mathbf{r}')$ goes to zero as $|\mathbf{r}-\mathbf{r}'|\rightarrow\infty$. The fundamental reason for this is the loss of quantum phase coherence between points that

are far apart. It is important to note that quantum locality also implies data locality, which is an essential property for scalable parallel algorithms. With parallel linear scaling methods, an n -fold increase in processors should lead approximately to an n -fold increase in simulation capability. In practice, however, this is true only for scalable algorithms. While there has been continued efforts to parallelize conventional quantum chemistry codes,^{33–41} taking advantage of quantum locality to achieve parallel efficiency is a largely unexplored but important area of research in linear scaling SCF theory.

One of most computationally demanding parts of a modern Gaussian-orbital SCF code is numerical integration of the exchange-correlation potential for DFT or HF/DFT calculations. Linear scaling approaches have been proposed to solve this problem.^{15–18} Most works use Becke⁴² nuclear weight functions that partition the integrand into a sum of single-center integrations, each with a nuclear origin. Standard numerical techniques may then be used to perform one-center numerical integration in spherical polar coordinates. It should be observed that this transformation overcomes the primary nuclear cusp in the electron density. However, there are additional complexities in the Becke approach. First, the

^{a)}Electronic mail: mchalla@lanl.gov

radial Jacobian cannot suppress cusps at other nuclear centers. These secondary cusps can be damped to zero with a sharper partition, but this has the risk of introducing additional discontinuities in the integrand. Second, increasing density of the radial grid does not guarantee convergence to a more accurate result. Third, the overlapping nature of the Becke weight functions impedes clean spatial partitioning, which is important for achieving true linear scaling, simple domain decomposition and data parallelism.

A different approach, hierarchical cubature (HiCu) employs an adaptive telescoping Cartesian grid with primitive cubature rules at the finest level of resolution.¹⁸ In HiCu, the hierarchical adaptive grid resolves strong variations and multiple length scales encountered in numerical integration of the exchange-correlation matrix \mathbf{K}_{xc} . The k -d tree^{43–45} data structure is used to represent the adaptive grid or CubeTree. The electron density of the system is also represented by the k -d tree data structure, called a RhoTree. Each node in a k -d tree contains a bounding box (BBox) that encloses the spatial data contained by it and its children (if the node is not a leaf node). Because of the hierarchical organization of data in a k -d tree, a search for all data within a given Euclidean distance (i.e., a range query) is an efficient $O(\log_2 N)$ operation.

The CubeTree is constructed through recursive bisection of a Cartesian domain (i.e., a box) enclosing the electron density, applying a fixed fully symmetric C_3 (cubature) rules⁴⁶ within each BBox. Given a local error threshold τ , the recursive bisection process stops when $\delta < \tau$, where

$$\delta = \Delta\rho + \frac{1}{3} \sqrt{(\Delta\rho_x)^2 + (\Delta\rho_y)^2 + (\Delta\rho_z)^2}. \quad (1)$$

Here $\Delta\rho$ is the magnitude of the difference between the exact electron charge in the BBox and the numerically integrated electron charge,

$$\Delta\rho = \left| \int_{\text{BBox}} \rho(\mathbf{r}) d\mathbf{r} - \sum_{i=1}^{N_g} w_i \rho(\mathbf{r}_i) \right|, \quad (2)$$

where w_i are the grid weights and N_g is the number of grid points in the cubature rule. The magnitude of the difference between the analytic and numerical integrations of the i th component (i.e., $i=x, y$ or z) of the density gradient $\nabla\rho$, denoted by $\Delta\rho_i$, is calculated in a fashion similar to $\Delta\rho$.

Unlike the first implementation of HiCu,¹⁸ the current implementation incorporates errors due to the density gradient, which leads to improved convergence properties for functionals that employ the generalized gradient approximation (GGA). Following the spirit of direct SCF methods,^{47,48} the accuracy of HiCu can be systematically improved by decreasing a threshold (τ). We also note that no fitting functions are employed. The nonoverlapping Cartesian approach combined with advanced data structures allows maximal exploitation of quantum locality, leading to an early onset of true linear scaling even for large basis sets and three-dimensional (3-D) systems. This nonoverlapping approach to quantum locality may be leveraged to achieve both data locality and straightforward domain decompositions.

Our purpose in this work is to describe our development of data parallel HiCu, with an emphasis on new, efficient and generally applicable load balancing algorithms. For the first

time, we propose measurement based load balancing (equal time partitioning) for domain decomposition in quantum chemistry. Equal time partitioning exploits the smooth variation of the electron density between SCF cycles (so called temporal locality⁴⁹), to predict an optimal domain decomposition for the next iteration.

The remainder of this paper is organized as follows: In Sec. II we discuss our strategies to efficiently parallelize the computation of \mathbf{K}_{xc} using HiCu. In Sec. III we discuss the issue of data locality to reduce communications between processors. In Sec. IV we describe a computational implementation of data parallel HiCu. In Sec. V we present the results and discussions of the speedup tests performed on three systems. In Sec. VI we summarize the main conclusions of the paper.

II. PARALLELIZATION OF HICU

In this section we describe the parallelization of HiCu with the goal of obtaining good speedups even for fine-grained parallelism, which we consider to be one heavy atom per processor. This is difficult, as the spatial distribution of the grid-work, which depends on the electron density, is highly irregular for an inhomogeneous system. We note that grid points can experience very different work loads depending on their environment, as the HiCu method will always try to minimize the work associated with gridding the density. To obtain a good speedup, we need to reduce both load imbalance and communications between processors. First we discuss the issue of load balancing. The issue of communication and data parallelism is discussed in Sec. III.

Since HiCu is purely Cartesian, spatial decomposition can be performed by simply sectioning the root BBox of the RhoTree into N_b nonoverlapping subboxes. The root BBox encloses the entire electron density, and is constructed by an upward merge of sub BBoxes in the RhoTree. At the lowest level, each BBox of the RhoTree is determined by the Gaussian extent.

One approach to load balance is the master-slave approach, which has been used to calculate \mathbf{K}_{xc} in conventional quantum chemistry programs.^{37,39,40} In the context of HiCu, a naive approach to master-slave load balancing involves making N_b larger than N_p (say, $N_b = 4N_p$), where N_p is the number of processors. The master program collects the results from a slave program which has finished its work on a particular subbox. The master program then instructs the slave program to work on a new subbox which has not been assigned to any of the slave programs before.

We have tested this simple dynamic master-slave load balancing scheme. Unfortunately, poor speedups are obtained, especially when both N_b and N_p are large. This is because as N_b increases, the time to deal with a subbox is generally small (since the volumes are smaller), so most slave programs will finish their work very quickly and they will spend significant time in contention for the master program, a well-known problem with this approach.^{50,51} In HiCu, this problem is exacerbated by the fact that the maximum load does not necessarily decrease linearly with N_b due to the irregular nature of the problem. It is interesting to note

that Guerra *et al.*³⁴ preferred a static dynamic load balancing scheme over master-slave load-balancing because they found that the repeated distribution of data requires much more communication.

An alternative is to employ a heuristic work estimate to guide the domain decomposition. We experimented with a number of heuristics, including splitting the density into equal-charge volumes and partitioning by center of mass. Of these heuristic approaches, the center of mass (CM) method, explained in the following, was found to have the greatest, but still unsatisfactory efficiency.

Another approach to load balance is measurement. Measurement approaches have several advantages: First they are easy to program. Second, they work well with smooth variation of the work load in the time or iteration space (temporal locality), an exploitable feature in the SCF procedure. Third, they have been employed with great success in tree-code methods for solving the N -body problem^{49,52–57}—HiCu has significant homology with the current implementation of the quantum chemical tree-code (QCTC)^{11,12,14} for solving the quantum Coulomb problem.

In the remainder of this work we describe our development of a measurement based method for load balancing data parallel HiCu. Briefly, the overall scheme is as follows: First the heuristic CM partitioning is used for the parallel calculation of \mathbf{K}_{xc} in the first SCF cycle. This provides a first measurement for repartitioning the root BBox in the next SCF cycle. As long as the starting guess is reasonable and pathological charge sloshing does not occur, each repartition should yield good load balance for the next iteration. In this way, temporal locality⁴⁹ of the problem is exploited to achieve a continuously improved load balance. Equal time repartitioning will be further explained in Sec. II B.

A. Center of mass partitioning scheme

In the center of mass (CM) partitioning scheme, we make a plausible assumption that the computing time for a subbox is proportional to the total electron charge in the box. To partition a BBox (see Fig. 1), we first calculate the center of the electron charge, which we have loosely called the “center of mass.” We assume that the electron charge due to an atom i is smeared out evenly in a sphere of radius R_i , which we have taken to be the Bragg–Slater radius⁵⁸ for the atom. The Bragg–Slater radius for an atom is the empirical atomic radius, derived from the fact that two atoms forming a bond in a crystal or molecule gives an approximate value of the internuclear distance. It is observed that interatomic bond lengths in solids and molecules, whether ionic or covalent, are approximately equal to pairwise sums of unique atomic radii.⁵⁸ The center of mass (charge) \mathbf{X} is given by

$$\mathbf{X} = \frac{\sum_i Z_{\text{eff},i} \mathbf{X}_{\text{eff},i}}{\sum_i Z_{\text{eff},i}}, \quad (3)$$

where $Z_{\text{eff},i} = Z_i(L_x L_y L_z / V_i)$, with $V_i = 4\pi R_i^3/3$. The symbols L_x , L_y , L_z and $\mathbf{X}_{\text{eff},i}$ are explained in Fig. 1. The index i in Eq. (3) runs through all atoms which overlap (totally or partially) with the box to be partitioned. In practice, we only calculate one component of \mathbf{X} , which is along the largest dimension of the box, since only one cutting plane is passing

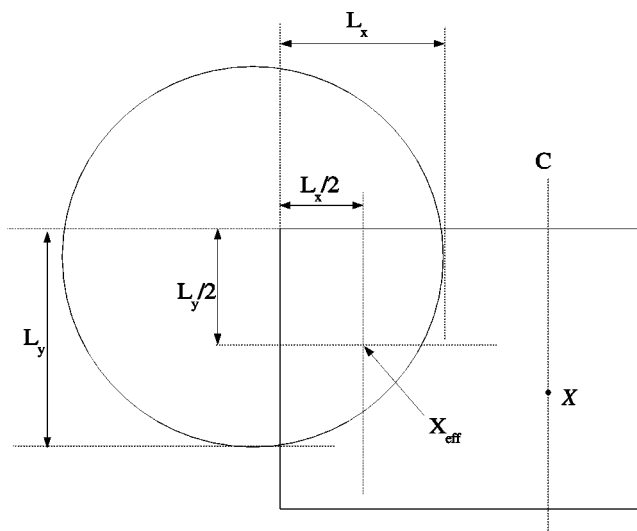


FIG. 1. A schematic diagram to explain the center of mass (CM) partitioning scheme. A sphere of Bragg–Slater (Ref. 58) radius R_i centered on the atom i is drawn for each atom (for simplicity, only one atom is considered in this case). The volume of the box–sphere intersection is approximated by $L_x L_y L_z$, where the L_x is the range of intersection between the sphere and the box in the x direction. L_y and L_z are calculated in a similar way. For simplicity we have assumed that the effective charge $Z_{\text{eff},i}$ is centered at \mathbf{X}_{eff} . Also shown is a cutting plane C passing through the center of mass \mathbf{X} .

through \mathbf{X} . Each of the two subboxes after a CM partitioning may be subjected to another CM partitioning. For simplicity, we always partition the root BBox in such a way that the number of the subboxes is a power of two, even though this restriction can be dropped easily (see Appendix A for a general equal time partitioning). In a parallel calculation, only one processor will perform the serial CM partition and broadcast the dimensions of subboxes to all other processors. It is important to emphasize that although the CM partitioning scheme might not give a good load balance, it does serve as a cheap and good starting partition for parallel HiCu. This is not a problem even for a large system because we can use a reasonably large local error threshold τ and a minimal basis set to start a calculation while the equal time partitioning scheme, to be explained in Sec. II B, will improve the efficiency in subsequent SCF cycles. One can then switch to a better basis set or lower the threshold τ during the iterative calculations.

B. Equal time partitioning scheme

Next, we discuss measurement approaches to domain decomposition, with the goal of repartitioning the root BBox based on workload information collected during CubeTree construction in the previous SCF cycle. By analogy with other parallel methods for computation of \mathbf{K}_{xc} , we first considered distributing an equal share of grid points, naively reasoning that the work load would be proportional to their number. If we evaluated the density using conventional methods, this would certainly be so. However, this is not a reliable tack because each leaf node interacts with the RhoTree in a different way, always attempting to minimize the total work load. Also, because each interaction of a leaf node with the RhoTree is so nonuniform, the simple interac-

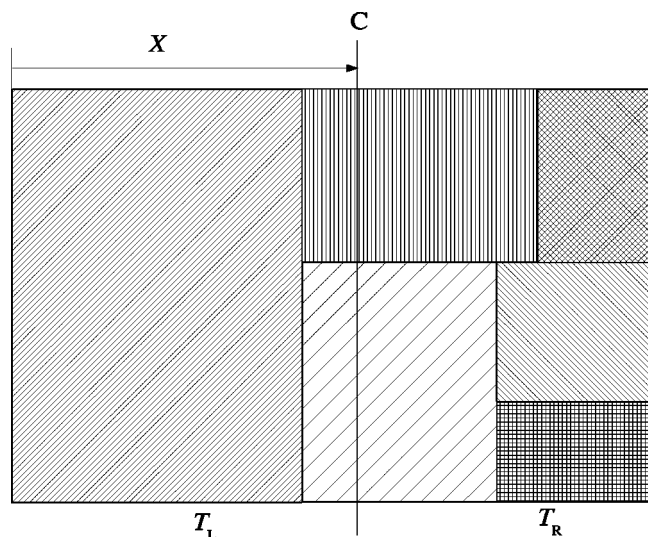


FIG. 2. A schematic diagram to illustrate the equal time (ET) partitioning scheme. For simplicity a two-dimensional analog is used. The rectangles represent the bounding “boxes” of the leaf nodes in a CubeTree. Different shades are used for different rectangles to indicate that the leaf-node times are different in general. The distance of the cutting line C from the left edge, x , is determined by the condition that the sum of the leaf-node times on both sides of the partition line C is equal (i.e., $T_L = T_R$). Notice that we have chosen the direction for partitioning as the direction which gives the largest dimension.

tion counting mechanisms employed by parallel N -body methods^{49,52–57} are not a viable alternative. This line of reasoning led us in the development of equal time (ET) partitioning, which uses the elapsed wall time associated with Cartesian sub-volumes (here the leaf nodes) to predict a good load balance in the next iteration.

In the ET partitioning of HiCu, the leaf-node time is stored along with other attributes in a variable of type CubeNode.¹⁸ With all processors holding leaf-node times in their respective local CubeTrees, ET creates a new domain decomposition by recursively partitioning a box into two subboxes such that each subbox carries approximately the same total leaf-node time (hence “equal time”). At the end of the decomposition we are left with N_p ET BBoxes in which to construct a HiCu grid. The astute reader will by now realize that this construction requires N_p to be a power of two, and that errors made early in the decomposition will propagate to the lowest level. The inconvenience associated with the restriction of a power of two for the number of processors is addressed in Appendix A. We have used a robust bisection method⁵⁹ to find the plane which equally (to some threshold) divides the workload in half (see Fig. 2). Notice that the cutting direction may be varied at each level of the recursion to avoid elongated ET BBoxes which could lead to poorly-balanced workloads. Our decomposition always chooses a direction which cuts along the largest box dimension, since the cubature rules will work best for leaf nodes close to cubic.⁴⁶

The application of ET to HiCu centers around interaction of the CubeTree leaf nodes with the RhoTree. This implementation does not account for the time taken to build each matrix element by traversing the CubeTree. This is because of the close homology between basis function products and

the electron density; they visit nearly identical portions of the CubeTree. However, if this homology did not exist, ET could associate primitive basis function products with Cartesian sub-volumes and wall times, using the combined information to load balance.

A working hypothesis of ET partitioning applied to HiCu is that the sum of leaf-node times for the root BBox is constant irrespective of sectioning. In practice this assumption is only approximately true, as the total number of leaf nodes and their times may fluctuate if the cutting planes are shifted. This is because with each new partitioning the CubeTree will readjust, potentially using a different number of leaf nodes to satisfy $\delta < \tau$. Also, as the time to create each leaf node depends on its size and position, readjustments will lead to leaf-node times that are somewhat different. Finally, variation of the density and hence the grid will occur in all but the last few SCF cycles. All of these factors work against ET partitioning. The magnitude of these factors depends on the granularity of the partition, as will be seen in Sec. V.

III. DATA PARALLEL APPROACH

Data parallelism is essential for scalable algorithms, as data replication stands to exhaust local memory and throttle communications with increasing N and N_p . For parallel HiCu, each of the ET BBoxes requires only local information about the electron density. By building a *locally essential RhoTree*, which is *just* sufficient to create the grid in an ET BBox, redundant replication of data is avoided. Another advantage is that the smaller locally essential RhoTree is more efficient to traverse relative to the entire RhoTree that would be encountered in a replicated data scheme. This reduced complexity due to data parallelism can lead to superlinear speedups as shown in Sec. V.

Data parallelism begins with a distributed atom-blocked compressed sparse row (DBC SR)⁶⁰ density matrix \mathbf{P}^p , which is used by program MAKERHO to create an intermediate distributed Hermite–Gaussian (HG)^{14,18,61} density ρ^p . The intermediate density consists of primitive distributions associated with only the rows spanned by \mathbf{P}^p . MAKERHO writes ρ^p to the local disk in parallel, avoiding IO and network bottlenecks. As an aside, this is the same density also used by program QCTC.^{11,12,14} The intermediate density ρ^p is generally insufficient to build a CubeTree in an ET BBox. Communication involving all processors now occurs to construct the locally essential density ρ_{LE}^p . Ultimately, sophisticated ordering and scheduling techniques could be used to reduce complexity of this step. So far, however, we have found no degradation of performance due to this redistribution. The construction of ρ_{LE}^p involves all processes going through their local primitive density distributions, performing box–box overlap tests as described in Ref. 18 to determine the exchange of data. After this redistribution of densities, the construction of the locally essential RhoTree, local CubeTree, and local partially summed exchange–correlation matrix (denoted by \mathbf{L}_{xc}^p) proceeds independently. The local, partially summed exchange–correlation matrix results from numerical integration of all primitive basis function products

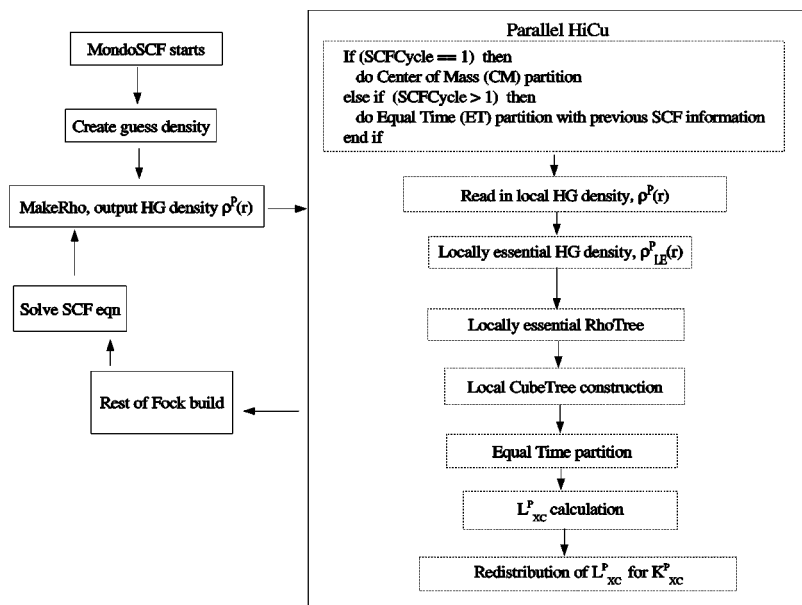


FIG. 3. The general flow of MONDOSCF, with the parallel HiCu subroutine expanded. Refer to the text for more details.

in the ET BBox, and is related to the full exchange-correlation matrix \mathbf{K}_{xc} by

$$\mathbf{K}_{xc} = \sum_{p=1}^{N_p} \mathbf{L}_{xc}^p, \quad (4)$$

although the summation in Eq. (4) is not explicitly carried out. Rather, \mathbf{L}_{xc}^p is redistributed and resumed to yield a non-overlapping row-distributed \mathbf{K}_{xc}^p .

As N_p becomes large, the ET BBox will become small and the corresponding local exchange-correlation matrix will become very sparse. During matrix construction there can be significant overheads associated with simple box–box testing of primitive basis function product BBox overlaps with the ET BBox. This problem is also encountered for periodic HiCu,⁶² where a double sum over lattice vectors leads to a substantial fraction of basis function products that are outside the root BBox. Our solution, described in detail elsewhere,⁶³ is to perform a high level test between basis function products in an atom–atom pair and the ET BBox using a cylinder-box overlap test.

Another potential complexity encountered with fine grained parallelism involves the redistribution and resummation of \mathbf{L}_{xc}^p . If one uses sparse-matrix data structures related to the standard CSR^{64,65} data structure such as the DBCSR or BCSR constructs,⁶⁰ the resummation of very many small sparse matrices can lead to large copy overheads. We have developed the Fast Matrix⁶⁶ (FastMat) data structure to address this problem, in which all rows are connected by a linear linked list, while each row is represented by a binary tree. This new construct has a low overhead for random insertions and updates of atom–atom blocks. It also allows many operations such as matrix–matrix addition and matrix truncation to be done in-place, reducing both memory and copy overheads.

IV. IMPLEMENTATION

We have implemented the parallel HiCu algorithm in MONDOSCF,⁶⁷ a suite of programs for linear scaling electronic structure theory and *ab initio* molecular dynamics. Figure 3 shows the general flow of an SCF cycle, with expanded detail of parallel HiCu. MONDOSCF has been written in Fortran 90/95 with the message-passing library MPI.⁶⁸ Leaf-node times are calculated with the MPI_WTIME MPI function.

V. RESULTS

We have performed scaling tests on taxol, a cluster of 110 water molecules, and α -pinene with 2^n processors, where n typically ranges from 1 to 7. These systems are chosen because they are highly inhomogeneous, three-dimensional systems posing a challenge to parallel HiCu. All runs on the first two systems were performed on a single SGI Origin 2000 with 128 (250 MHz MIPS R10000 MIPS) processors. For the calculations on α -pinene, we used a cluster of 256 4 CPU HP/Compaq Alphaservert ES45s with the Quadrics QsNet High Speed Interconnect.

We start the calculations with the STO-3G basis set and a large threshold τ , switch to the 3-21G basis set using our mixed integral approach, and run for three SCF cycles. The density matrix \mathbf{P} is saved to disk and scaling tests of parallel HiCu are performed. Timings for parallel HiCu are taken at the fourth cycle where the center of mass partition is used, and at the fifth cycle where the equal time partition is used. The procedure of switching the basis set and iterating for a few SCF cycles places the electron density in the middle of convergence. This is not really necessary, as performance is quite insensitive to the SCF cycle.

The result of taxol scaling tests is shown in Fig. 4, where the center of mass partitions perform rather poorly past 16 processors, with a speedup of only 22.4 obtained at 64 processors. However, the equal time partitions give good scal-

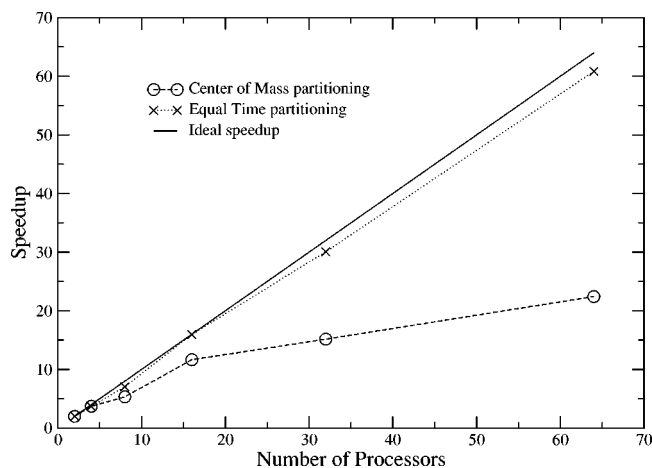


FIG. 4. Scaling of parallel HiCu on taxol ($C_{47}H_{51}NO_{14}$) RPBE/3-21G, with $\tau = 1.0 \times 10^{-6}$. Speedups are relative to a 2-processor calculation.

ing results compared to the ideal speedup, obtaining a speedup of 60.8 with 64 processors. We note that our result of a speedup of 7.03 with 8 processors compares favorably with the speedup of about 6.0 of Sosa *et al.*,³⁸ which is for an entire single-point energy calculation.

Similar scaling tests have been performed on a cluster of 110 water molecules. Figure 5 shows that the center of mass partition gives a better speedup than in the case of taxol. Equal time repartitioning is highly effective. The speedup is almost ideal up to 64 processors. We observe a superlinear speedup at 32 processors, which is a medium-grained calculation. We attribute the superlinear speedup to a good load balance and to reduced complexity associated with data parallelism, where the memory access time has been greatly reduced with a smaller RhoTree. The speedup with 128 processors degrades slightly to 112.64 (an efficiency of 88%). Still, this is an encouraging result as the 128-processor run corresponds to very fine-grained parallelism, with less than one heavy atom per processor. The reason for this degraded efficiency is not due to communication overhead but to intrinsic limitations in the ET scheme associated with fluctuations in the leaf-node count and small nonconservation of

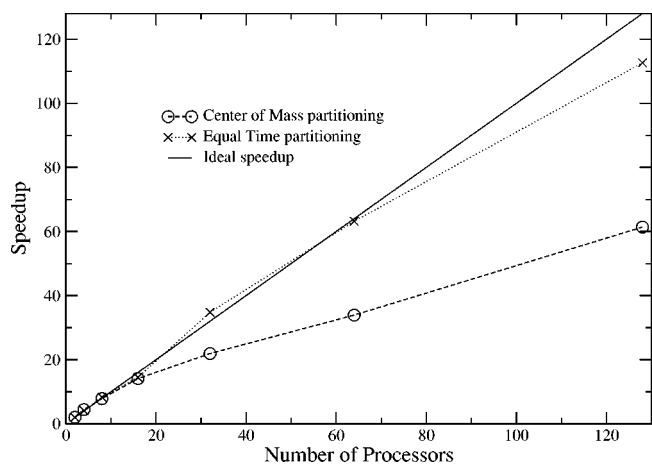


FIG. 5. Scaling of parallel HiCu on $(H_2O)_{110}$ RPBE/3-21G, with $\tau = 1.0 \times 10^{-6}$. Speedups are relative to a 2-processor calculation.

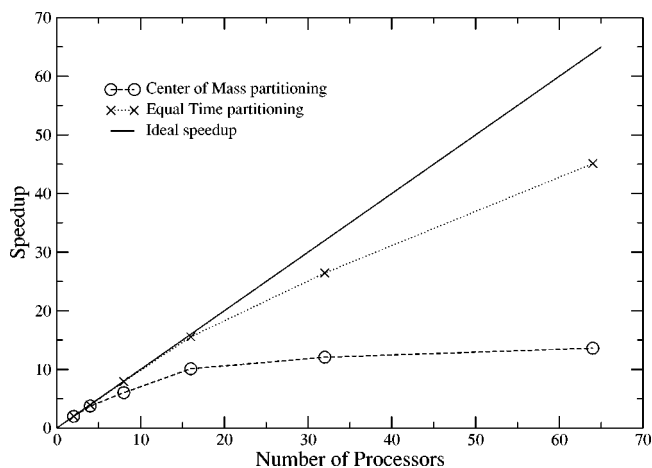


FIG. 6. Scaling of parallel HiCu on α -pinene RPBE/6-311G(df,p), with $\tau = 1.0 \times 10^{-6}$. Speedups are relative to a 2-processor calculation.

total times, as explained in Sec. II B. While there are almost certainly ways to overcome/compensate for these effects, we feel that the overall performance of ET for parallel HiCu is outstanding. To explore the behavior of the current implementation with $N_p \sim 1024$ will require a fully parallel MONDOSCF, which we are actively pursuing.

A well studied system in parallel quantum chemistry is 6-311G(df,p) α -pinene ($C_{10}H_{16}$).^{37,38} To date, the most competitive rate for this system is 37.3 out of 64 processors, achieved by Furlani *et al.*³⁷ over a complete SCF cycle. Scaling results for a medium accuracy grid corresponding to $\tau = 1.0 \times 10^{-6}$ are shown in Fig. 6, with HiCu achieving a competitive rate of 45.11 (70% efficiency) out of 64 processors. More impressively, with a high accuracy grid ($\tau = 1.0 \times 10^{-8}$) equal time partitioning delivers an unrivaled 99.5% efficiency with $N_p = 64$ as shown in Fig. 7. These results demonstrate that the efficiency of ET improves when the grid work per processor is increased. This is not surprising as increasing the grid work decreases the effective granularity, reducing sensitivity to fluctuations in leaf-node times and leading to a more predictive equal time partition.

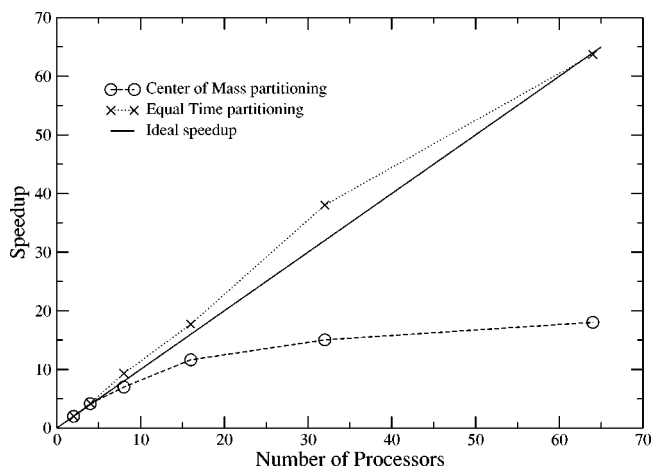


FIG. 7. Scaling of parallel HiCu on α -pinene RPBE/6-311G(df,p), with $\tau = 1.0 \times 10^{-8}$. Speedups are relative to a 2-processor calculation.

VI. CONCLUSIONS

We have developed an efficient and general load balancing technique, equal time (ET) partitioning, and used it in data parallel implementation of the linear scaling hierarchical cubature (HiCu) algorithm for numerical integration of the exchange-correlation matrix. In the ET approach, strong spatial irregularities in the work load are overcome by exploiting temporal locality between iterations. We expect ET to also exploit this effect between geometry steps in an optimization or molecular dynamics run. In this data parallel implementation, quantum locality has been used to reduce communications between processors, to lower memory usage and to reduce the computational complexity associated with traversal of the RhoTree. In some cases, this last effect can lead to superlinear speedups for medium grained calculations.

To the best of our knowledge, equal time partitioning is the first measurement based load balancing strategy applied to quantum chemistry. While ET partitioning shares some commonalities with partitioning schemes for parallel N -body tree-code solvers such as orthogonal recursive bisection (ORB)⁵² and Costzones,⁵⁵ ET differs by (i) using exact load timing information rather than counting the number of interactions as in the ORB and Costzones methods, and (ii) associating this load with rectilinear partitioning of Cartesian space rather than particle groupings (ORB) or nodes in a tree (Costzones). These are important distinctions for quantum chemistry. Often methods are much more complicated than the simple N -body problem. Measured loads encompass irregularities associated with different Gaussian extents, non-uniform access patterns, etc. that cannot be accounted for by counting interactions. Also, because we are often doing several complicated things at once, attaching work loads to a portion of 3-D space is a convenient catch-all that can be associated with most objects and processes: basis function products and matrix construction, electron densities and grid construction, etc.

Our parallel approach to construction of \mathbf{K}_{xc} differs fundamentally from other all-electron approaches in that (i) we do not employ Becke weights and (ii) we employ measurement based load balancing. It is also different from many implementations in that (iii) it is data parallel. In addition to the numerical issues discussed in the Introduction, we believe the nonoverlapping Cartesian subdivision used by HiCu enables a more straightforward approach to data locality and domain decomposition. Both the replicated data and the master-slave approach have well known shortcomings for fine grained massive parallelism, which have been avoided here.

With $N_p \leq 128$, the predominant cause of deviation from ideal scaling is due to basic nonlinearities in ET partitioning of the HiCu grid, rather than overheads associated with inefficient data structures or communication. These deviations increase as the processor to grid work ratio increases, i.e., with finer granularity. The main factors that determine grid work are the number of atoms and τ , the threshold controlling grid accuracy. For example, calculations on $(\text{H}_2\text{O})_{110}$ with a medium accuracy grid and ~ 1 heavy atom per processor are fine grained, achieving a parallel efficiency of 88% with $N_p = 128$. Even with $N_p/N_{\text{Atoms}} > 2$ it is possible

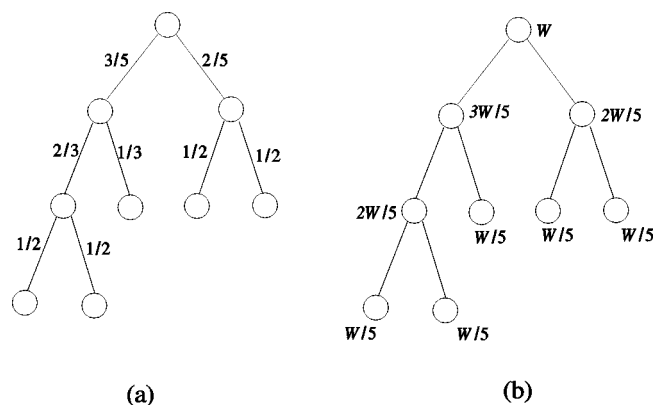


FIG. 8. An illustrative example to show the general equal time partitioning scheme for the case of $m=5$.

achieve excellent scaling by lowering τ . This was demonstrated by 6-311G(df,p) calculations on α -pinene ($\text{C}_{10}\text{H}_{16}$) with $N_p = 64$, achieving a parallel efficiency of 99.5%. Apparently, this is the first time this effect has been demonstrated in the context of parallel quantum chemistry. More importantly, because HiCu continues to achieve early onset linear scaling even with tight numerical thresholds,¹⁸ these results suggest the potential for very large, high quality calculations through a combination of massive parallelism and scalable $O(N)$ algorithms.

ACKNOWLEDGMENTS

This work has been carried out under the auspices of the U.S. Department of Energy under Contract No. W-7405-ENG-36 and the ASCI project. The Advanced Computing Laboratory of Los Alamos National Laboratory, Los Alamos, NM 87545 is acknowledged. This work was performed on computing resources located at this facility. We thank the reviewer for valuable comments. A fruitful discussion on the general equal time partitioning scheme with C. J. Tymczak is appreciated.

APPENDIX: A GENERAL EQUAL TIME PARTITIONING SCHEME

In this Appendix a general equal time partitioning scheme is described where the number of subboxes, denoted by m , is not restricted to a power of two. To facilitate this general partitioning, first create a binary tree with m leaf nodes, with the condition that the depth between any two leaf nodes differs by no more than one. Next, associate the edges connecting a parent node to its left and right child nodes with edge weights equal to $n_L/(n_L+n_R)$ and $n_R/(n_L+n_R)$, respectively. Here n_L (n_R) is the number of left (right) leaf nodes of the parent node. Figure 8(a) shows a binary tree and all the edge weights for the case of $m=5$. With these edge weights, it is a fairly simple matter to partition a total workload W into m equal workloads. Starting from the root node, recursively divide the workload in a parent node according to the ratio of the edge weights, i.e., $n_L:n_R$. At the end of this procedure, all leaf nodes will have a workload of W/m . Fig-

ure 8(b) shows the detailed workload information in all nodes for the case of $m = 5$. In a domain decomposition of a root BBox, a bisection search may be used to partition a box into two subboxes with a workload ratio equal to $n_L:n_R$, which can be easily read off from the edge weights of the binary tree. It should be noticed that the “power-of-two” equal time partitioning scheme explained in Sec. II B is a special case of this general partitioning scheme.

- ¹S. Goedecker, *Rev. Mod. Phys.* **71**, 1085 (1999).
- ²S. Y. Wu and C. S. Jayanthi, *Phys. Rep.* **358**, 1 (2002).
- ³E. Schwegler and M. Challacombe, *J. Chem. Phys.* **105**, 2726 (1996).
- ⁴E. Schwegler, M. Challacombe, and M. Head-Gordon, *J. Chem. Phys.* **106**, 9708 (1997).
- ⁵E. Schwegler, M. Challacombe, and M. Head-Gordon, *J. Chem. Phys.* **109**, 8764 (1998).
- ⁶E. Schwegler, “Thesis: Linear scaling computation of the Hartree–Fock exchange matrix,” Ph.D. thesis, University of Minnesota, 1998, URL <http://www.t12.lanl.gov/~mchalla/>
- ⁷E. Schwegler and M. Challacombe, *J. Chem. Phys.* **111**, 6223 (1999).
- ⁸E. Schwegler and M. Challacombe, *Theor. Chem. Acc.* **104**, 344 (2000).
- ⁹C. A. White, B. Johnson, P. Gill, and M. Head-Gordon, *Chem. Phys. Lett.* **230**, 8 (1994).
- ¹⁰C. A. White, B. G. Johnson, P. M. W. Gill, and M. Head-Gordon, *Chem. Phys. Lett.* **253**, 268 (1996).
- ¹¹M. Challacombe, E. Schwegler, and J. Almlöf, *Computational Chemistry: Review of Current Trends* (World Scientific, Singapore, 1996), pp. 53–107.
- ¹²M. Challacombe, E. Schwegler, and J. Almlöf, *J. Chem. Phys.* **104**, 4685 (1996).
- ¹³M. C. Strain, G. E. Scuseria, and M. J. Frisch, *Science* **271**, 51 (1996).
- ¹⁴M. Challacombe and E. Schwegler, *J. Chem. Phys.* **106**, 5526 (1997).
- ¹⁵J. M. Pérez-Jordá and W. Yang, *Chem. Phys. Lett.* **241**, 469 (1995).
- ¹⁶R. E. Stratmann, G. E. Scuseria, and M. J. Frisch, *Chem. Phys. Lett.* **257**, 213 (1996).
- ¹⁷C. F. Guerra, J. G. Snijders, G. te Velde, and E. J. Baerends, *Theor. Chem. Acc.* **99**, 391 (1998).
- ¹⁸M. Challacombe, *J. Chem. Phys.* **113**, 10037 (2000).
- ¹⁹X. P. Li, R. W. Nunes, and D. Vanderbilt, *Phys. Rev. B* **47**, 10891 (1993).
- ²⁰M. S. Daw, *Phys. Rev. B* **47**, 10895 (1993).
- ²¹S. Y. Qiu, C. Z. Wang, K. M. Ho, and C. T. Chan, *J. Phys.: Condens. Matter* **6**, 9153 (1994).
- ²²E. Hernández and M. J. Gillan, *Phys. Rev. B* **51**, 10157 (1995).
- ²³E. Hernández, M. J. Gillan, and C. Goringe, *Phys. Rev. B* **53**, 7147 (1996).
- ²⁴C. M. Goringe, E. Hernández, M. J. Gillan, and I. J. Bush, *Comput. Phys. Commun.* **94**, 89 (1996).
- ²⁵A. D. Daniels, J. M. Millam, and G. E. Scuseria, *J. Chem. Phys.* **107**, 425 (1997).
- ²⁶D. R. Bowler and M. J. Gillan, *Comput. Phys. Commun.* **120**, 95 (1999).
- ²⁷A. H. R. Palser and D. E. Manolopoulos, *Phys. Rev. B* **58**, 12704 (1998).
- ²⁸M. Challacombe, *J. Chem. Phys.* **110**, 2332 (1999).
- ²⁹A. M. N. Niklasson, *Phys. Rev. B* **66**, 155115 (2002).
- ³⁰A. M. N. Niklasson, C. J. Tymczak, and M. Challacombe, *J. Chem. Phys.* (to be published).
- ³¹W. Kohn, *Int. J. Quantum Chem.* **56**, 229 (1995).
- ³²W. Kohn, *Phys. Rev. Lett.* **76**, 3168 (1996).
- ³³R. J. Harrison and R. Shepard, *Annu. Rev. Phys. Chem.* **45**, 623 (1994).
- ³⁴C. F. Guerra, O. Visser, J. G. Snijders, G. te Velde, and E. J. Baerends, in *Methods and Techniques for Computational Chemistry*, edited by E. Clementi and G. Corongiu (STEF, Cagliari, 1995), p. 305.
- ³⁵C. P. Sosa, J. Ochterski, J. Carpenter, and M. J. Frisch, *J. Comput. Chem.* **19**, 1053 (1998).
- ³⁶F. Stephan and W. Wenzel, *J. Chem. Phys.* **108**, 1015 (1998).
- ³⁷T. R. Furlani, J. Kong, and P. M. W. Gill, *Comput. Phys. Commun.* **128**, 170 (2000).
- ³⁸C. P. Sosa, G. Scalmani, R. Gomperts, and M. J. Frisch, *Parallel Comput.* **26**, 843 (2000).
- ³⁹T. Yoshihiro, F. Sato, and H. Kashiwagi, *Chem. Phys. Lett.* **346**, 313 (2001).
- ⁴⁰J. Baker and P. Pulay, *J. Comput. Chem.* **23**, 1150 (2002).
- ⁴¹H. Takashima, S. Yamada, S. Obara, K. Kitamura, S. Inabata, N. Miyakawa, K. Tanabe, and U. Nagashima, *J. Comput. Chem.* **23**, 1337 (2002).
- ⁴²A. D. Becke, *J. Chem. Phys.* **88**, 2547 (1988).
- ⁴³J. L. Bentley and J. H. Friedman, *ACM Comput. Surv.* **11**, 397 (1979).
- ⁴⁴J. L. Bentley, *Commun. ACM* **23**, 214 (1980).
- ⁴⁵V. Gaede and O. Günther, *ACM Comput. Surv.* **30**, 170 (1998).
- ⁴⁶A. H. Stroud, *Approximate Calculation of Multiple Integrals*, Prentice-Hall series in automatic computation (Prentice-Hall, Englewood Cliffs, NJ, 1971).
- ⁴⁷J. Almlöf, K. Faegri, and K. Korsell, *J. Comput. Chem.* **3**, 385 (1982).
- ⁴⁸M. Häser and R. Ahlrichs, *J. Comput. Chem.* **10**, 104 (1989).
- ⁴⁹J. R. Pilkington and S. B. Baden, *IEEE Trans. Para. Distr. Sys.* **7**, 288 (1996).
- ⁵⁰B. Wilkinson and M. Allen, *Parallel Programming* (Prentice-Hall, Upper Saddle River, NJ, 1999).
- ⁵¹G. V. Wilson, *Practical Parallel Programming* (The MIT Press, Cambridge, MA, 1995).
- ⁵²M. S. Warren and J. K. Salmon, *Proceedings of Supercomputing '92*, 1992, p. 570.
- ⁵³A. Y. Grama, V. Kumar, and A. Sameh, *Proceedings of Supercomputing '94*, 1994, p. 439.
- ⁵⁴M. S. Warren and J. K. Salmon, *Comput. Phys. Commun.* **87**, 266 (1995).
- ⁵⁵J. P. Singh, C. Holt, J. L. Hennessy, and A. Gupta, *Proceedings of Supercomputing '93*, 1993, p. 54.
- ⁵⁶J. P. Singh, C. Holt, T. Totsuka, A. Gupta, and J. Hennessy, *J. Para. Distr. Comput.* **27**, 118 (1995).
- ⁵⁷A. Grama, V. Kumar, and A. Sameh, *Parallel Comput.* **24**, 797 (1998).
- ⁵⁸J. C. Slater, *J. Chem. Phys.* **41**, 3199 (1964).
- ⁵⁹W. H. Press, S. A. Teukolsky, W. T. Vetterling, and B. P. Flannery, *Numerical Recipes in FORTRAN* (Cambridge University Press, Port Chester, NY, 1992).
- ⁶⁰M. Challacombe, *Comput. Phys. Commun.* **128**, 93 (2000).
- ⁶¹G. R. Ahmadi and J. Almlöf, *Chem. Phys. Lett.* **246**, 364 (1995).
- ⁶²C. J. Tymczak and M. Challacombe, “Linear scaling computation of the Fock matrix. VII. Periodic boundary conditions,” in preparation.
- ⁶³M. Challacombe and C. J. Tymczak, “Data structures and error estimates for linear scaling SCF theory. I. Trees, potentials and the Gaussian extent,” in preparation.
- ⁶⁴F. G. Gustavson, *ACM Trans. Math. Softw.* **4**, 250 (1978).
- ⁶⁵S. Pissanetzky, *Sparse Matrix Technology* (Academic, London, 1984).
- ⁶⁶M. Challacombe, C. K. Gan, and A. M. N. Niklasson, “Data structures and error estimates for linear scaling SCF theory. II. Sparse matrices, atom-blocking, truncation and inverse Cholesky factors,” in preparation.
- ⁶⁷M. Challacombe, E. Schwegler, C. Tymczak, C. K. Gan, K. Nemeth, A. M. N. N. H. Nymeyer, and G. Henkleman, MONDOSCF v1.0a7, “A program suite for massively parallel, linear scaling SCF theory and *ab initio* molecular dynamics,” 2001, URL <http://www.t12.lanl.gov/~mchalla/>, Los Alamos National Laboratory (LA-CC 01-2), Copyright University of California.
- ⁶⁸“Message Passing Interface Forum. MPI: A message-passing interface standard (version 2.0),” 1998, <http://www.mpi-forum.org>