

# A simplified density matrix minimization for linear scaling self-consistent field theory

Matt Challacombe

Citation: *The Journal of Chemical Physics* **110**, 2332 (1999); doi: 10.1063/1.477969

View online: <https://doi.org/10.1063/1.477969>

View Table of Contents: <http://aip.scitation.org/toc/jcp/110/5>

Published by the [American Institute of Physics](#)

---

## Articles you may be interested in

[Linear scaling conjugate gradient density matrix search as an alternative to diagonalization for first principles electronic structure calculations](#)

*The Journal of Chemical Physics* **106**, 5569 (1997); 10.1063/1.473579

[Trace resetting density matrix purification in  \$O\(N\)\$  self-consistent-field theory](#)

*The Journal of Chemical Physics* **118**, 8611 (2003); 10.1063/1.1559913

[Communication: Generalized canonical purification for density matrix minimization](#)

*The Journal of Chemical Physics* **144**, 091102 (2016); 10.1063/1.4943213

[Semiempirical methods with conjugate gradient density matrix search to replace diagonalization for molecular systems containing thousands of atoms](#)

*The Journal of Chemical Physics* **107**, 425 (1997); 10.1063/1.474404

[Linear and sublinear scaling formation of Hartree–Fock-type exchange matrices](#)

*The Journal of Chemical Physics* **109**, 1663 (1998); 10.1063/1.476741

[What is the best alternative to diagonalization of the Hamiltonian in large scale semiempirical calculations?](#)

*The Journal of Chemical Physics* **110**, 1321 (1999); 10.1063/1.478008

---

PHYSICS TODAY

WHITEPAPERS

### ADVANCED LIGHT CURE ADHESIVES

Take a closer look at what these environmentally friendly adhesive systems can do

READ NOW

PRESENTED BY  
 **MASTERBOND**  
ADHESIVES | SEALANTS | COATINGS

# A simplified density matrix minimization for linear scaling self-consistent field theory

Matt Challacombe<sup>a)</sup>

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

(Received 1 September 1998; accepted 7 October 1998)

A simplified version of the Li, Nunes and Vanderbilt [Phys. Rev. B **47**, 10891 (1993)] and Daw [Phys. Rev. B **47**, 10895 (1993)] density matrix minimization is introduced that requires four fewer matrix multiplies per minimization step relative to previous formulations. The simplified method also exhibits superior convergence properties, such that the bulk of the work may be shifted to the quadratically convergent McWeeny purification, which brings the density matrix to idempotency. Both orthogonal and nonorthogonal versions are derived. The AINV algorithm of Benzi, Meyer, and Tuma [SIAM J. Sci. Comp. **17**, 1135 (1996)] is introduced to linear scaling electronic structure theory, and found to be essential in transformations between orthogonal and nonorthogonal representations. These methods have been developed with an atom-blocked sparse matrix algebra that achieves sustained megafloating point operations per second rates as high as 50% of theoretical, and implemented in the MondoSCF suite of linear scaling SCF programs. For the first time, linear scaling Hartree–Fock theory is demonstrated with three-dimensional systems, including water clusters and estane polymers. The nonorthogonal minimization is shown to be uncompetitive with minimization in an orthonormal representation. An early onset of linear scaling is found for both minimal and double zeta basis sets, and crossovers with a highly optimized eigensolver are achieved. Calculations with up to 6000 basis functions are reported. The scaling of errors with system size is investigated for various levels of approximation. © 1999 American Institute of Physics. [S0021-9606(99)30702-9]

## I. INTRODUCTION

Computation of the Fock matrix has historically been the limiting step in quantum chemical applications of the Hartree–Fock (HF) and Kohn–Sham (KS) theories to large systems. This is due to the expensive  $\mathcal{O}(N_{\text{bas}}^2)$  cost of two-electron integral computation and manipulation<sup>1</sup> in the direct method,<sup>2,3</sup> where  $N_{\text{bas}}$  is the number of basis functions and is proportional to system size. Recently, methods with a computational complexity of  $\mathcal{O}(N_{\text{bas}})$  have been introduced for computing the Fock matrix,<sup>4–16</sup> allowing access to systems large enough that the  $\mathcal{O}(N_{\text{bas}}^3)$  eigensolution of the self-consistent field (SCF) equations is presently the bottleneck in the large scale application of HF and KS theories.

The computationally demanding aspects of quantum chemical SCF theory are now identical to those encountered in condensed matter formulations of the local density approximation (LDA) and tight binding theories, for which a variety of linear scaling algorithms exist. These include density matrix minimization (DMM),<sup>17–24</sup> orbital minimization (OM),<sup>25–29</sup> and the Fermi operator expansion (FOE).<sup>30–35</sup> See Refs. 36–38 for excellent reviews and comparisons of these methods, as well as a more complete bibliography.

These methods achieve linear scaling for insulating systems and finite temperature metals by exploiting the short range nature of quantum interactions, the locality of which

manifests itself in exponential decay of the density matrix.<sup>38–40</sup>

$$\rho(\mathbf{r}, \mathbf{r}') \sim \exp(-\sqrt{E_{\text{gap}}}|\mathbf{r} - \mathbf{r}'|), \quad (1)$$

where  $E_{\text{gap}}$  is the HOMO-LUMO gap, which is the difference between eigenvalues of the highest occupied molecular orbital and the lowest unoccupied molecular orbital. If a basis of local functions  $\phi$  is used in calculation of the density matrix.

$$\rho(\mathbf{r}, \mathbf{r}') = \sum_{ij} P_{ij} \phi_i(\mathbf{r}) \phi_j(\mathbf{r}'), \quad (2)$$

then it follows that the discrete representation  $\mathbf{P}$  has similar localization properties. This has been observed for both HF<sup>14,40</sup> and KS<sup>38</sup> models. Thus, for a sufficiently large system only  $\mathcal{O}(N_{\text{bas}})$  matrix elements are expected to remain numerically significant, and sparse matrix methods may be used to obtain the SCF in linear scaling CPU time.

The DMM, OM, and FOE methods employ a single variable, such as a localization radius or a matrix element threshold, which defines an approximate sparse matrix algebra, and which allows deviation from the exact result to be controlled systematically. However, the extension of these methods to the quantum chemical domain has faced a number of challenges related to the use of large nonorthogonal basis sets, an exacting demand for error control, and the slow (approximate)

<sup>a)</sup>Electronic mail: MChalla@T12.LANL.Gov

mately exponential but slow) decay in off-diagonal elements of  $\mathbf{P}$  with increasing system size,<sup>14,40</sup> which can be severe in the case of three dimensions.

Clearly, the off-diagonal behavior of  $\mathbf{P}$  depends on the physical system, the model chemistry, and the basis set. For example, minimal basis set calculations may lead to large HOMO-LUMO gaps. Likewise, HF theory tends to overestimate the HOMO-LUMO gap while LDA tends to underestimate it. One approach to limiting the bandwidth of  $\mathbf{P}$  is to construct basis sets with only local support.<sup>41,42</sup> Another approach is the use of orthogonal bases such as wavelets<sup>43-45</sup> or grids<sup>46-48</sup> which tend to yield very sparse matrices. A less exotic approach, pursued here, is to use standard quantum chemical basis sets together with efficient numerical algorithms to access system sizes that yield a sparse representation.

To this end, a simplified version of the Li, Nunes, Vanderbilt<sup>17</sup> and Daw<sup>18</sup> (LNVD) density matrix minimization (SDMM) is introduced that requires four fewer matrix multiplications than standard versions of the density matrix minimization, and which permits the time dominant work to be carried out with the quadratically convergent McWeeny purification.<sup>49</sup> The simplified density matrix minimization is implemented in both an orthogonal and nonorthogonal representation using an atom-blocked sparse matrix algebra that exploits both data locality and available sparsity. The AINV<sup>50,51</sup> method for computation of the incomplete inverse Cholesky factor is introduced to electronic structure theory as an efficient and accurate method for enabling the congruence transformation to and from an orthogonal representation, and for applying the inverse overlap matrix in nonorthogonal methods.

The paper is organized as follows. In Sec. II, the conventional and LNVD approach to solving the SCF equations is reviewed. In Sec. III, details of the LNVD method are presented, and a diagonal guess is shown to greatly simplify gradient and line searches. In Sec. IV, the sparse atom-blocked matrix algebra is introduced, and details of its implementation in the MondoSCF suite of linear scaling SCF programs<sup>52</sup> are given. Then in Sec. V, timings and errors are presented and discussed for water clusters and estane polymers. Finally, conclusions are drawn in Sec. VI.

## II. OVERVIEW

### A. SCF theory

In SCF theory, the density matrix  $\mathbf{P}$  defines the entire physical system.<sup>53</sup> In particular, the electronic energy<sup>54</sup> is

$$E_{\text{el}} = \text{Tr}\{(\mathbf{h} + \mathbf{F}[\mathbf{P}])\mathbf{P}\}, \quad (3)$$

where  $\mathbf{F}$  is the Fock matrix, which depends on  $\mathbf{P}$  in a complicated way, and  $\mathbf{h}$  is the core Hamiltonian, which does not depend on  $\mathbf{P}$ . For mathematical convenience, the closed shell density matrix may be factored in terms of the occupied molecular orbital coefficients  $\mathbf{C}_{\text{occ}}$  as

$$\mathbf{P} = 2\mathbf{C}_{\text{occ}}\mathbf{C}_{\text{occ}}^\dagger. \quad (4)$$

The molecular orbital matrix  $\mathbf{C}$ , including both occupied and virtual orbitals, is orthogonal;

$$\mathbf{C}^\dagger\mathbf{C} = \mathbf{I}. \quad (5)$$

The minimization of  $E_{\text{el}}$  with respect to  $\mathbf{C}$ , and with the constraint of orthogonality, leads to Roothaan's form of the SCF equations<sup>55,54</sup>

$$\mathbf{F}[\mathbf{P}^{n-1}]\mathbf{C}^n = \mathbf{C}^n\epsilon, \quad (6)$$

which is an eigenproblem with  $\mathcal{O}(N_{\text{bas}}^3)$  computational complexity. Iteration of Eq. (6) leads to the minimum energy solution at self-consistency, when  $\mathbf{P}^n = \mathbf{P}^{n-1}$ . For well behaved problems, typically  $n \sim 10$  SCF cycles are required to reach convergence.

In the preceding matrix equations, an orthogonal basis has been tacitly assumed. In practice however, matrix elements are evaluated in the nonorthogonal basis of atomic orbitals (AOs). Introducing the metric or overlap matrix with elements

$$S_{ab} = (\phi_a, \phi_b) \quad (7)$$

between AO basis functions  $\phi_a$  and  $\phi_b$ , the Roothaan equations may be posed as the generalized eigenvalue problem<sup>55</sup>

$$\mathbf{F}_{\text{ao}}\mathbf{C}_{\text{ao}} = \mathbf{S}_{\text{ao}}\epsilon, \quad (8)$$

where subscript ao implies a nonorthogonal representation.

### B. Orthogonal and nonorthogonal bases

It is well known that the generalized eigenproblem, Eq. (8), can be transformed to the standard eigenproblem, Eq. (6), using a congruence transformation involving factorization of the metric (overlap) matrix.<sup>56,57</sup> A key matrix in this transformation is  $\mathbf{Z}$ , the inverse factor, which relates nonorthogonal (AO) and orthogonal representations of the Fock matrix

$$\mathbf{F} = \mathbf{Z}\mathbf{F}_{\text{ao}}\mathbf{Z}^T \quad (9)$$

and the density matrix

$$\mathbf{P}_{\text{ao}} = \mathbf{Z}^T\mathbf{P}\mathbf{Z}. \quad (10)$$

The inverse factor has the property

$$\mathbf{Z}^T\mathbf{S}\mathbf{Z} = \mathbf{I}, \quad (11)$$

which is satisfied by  $\mathbf{Z} = \mathbf{S}^{-1/2}$  or  $\mathbf{Z} = \mathbf{L}^{-1}$ , where  $\mathbf{L}$  is the Cholesky factor for which  $\mathbf{S} = \mathbf{L}\mathbf{L}^T$ . The choice  $\mathbf{Z} = \mathbf{S}^{-1/2}$  corresponds to Löwdin's symmetric orthogonalization,<sup>58,59</sup> which is commonly used in modern quantum chemistry. The choice  $\mathbf{Z} = \mathbf{L}^{-1}$  is widely used in solution of the generalized eigenproblem,<sup>56,57</sup> and has recently been introduced to linear scaling SCF theory by Millam and Scuseria (MS).<sup>23</sup>

Symmetric orthogonalization requires an eigensolve, which is  $\mathcal{O}(N^3)$ . As discussed by MS,<sup>23</sup> the incomplete Cholesky factorization is potentially  $\mathcal{O}(N)$  when sparse matrix methods are employed in conjunction with an elimination tree. In the MS approach, the approximate inverse factor  $\mathbf{Z}$  is obtained through an incomplete linear solve. This tends to be inaccurate unless tight thresholds are employed, in which case the solve may become expensive. Moreover, inverting the Cholesky factor introduces two levels of incompleteness (one in the factorization and one in the inversion), which may introduce errors that are difficult to control.

A more elegant approach, pursued here, is to solve for  $\mathbf{Z}=\mathbf{L}^{-1}$  directly using the AINV algorithm.<sup>50,51</sup> AINV is state of the art in the theory and application of inverse preconditioners, and may be used to obtain the inverse factors  $\mathbf{Z}$  to within arbitrary accuracy. The matrix  $\mathbf{Z}$  is obtained by S orthogonalization of the standard basis vector, where sparsity is preserved with the use of a drop tolerance.<sup>50,51</sup> We have recently implemented a block version of AINV<sup>60</sup> and found it to be very fast under the demands of high accuracy.

When working in a nonorthogonal representation, the inverse overlap matrix  $\mathbf{S}^{-1}$  enters the OM,<sup>25,26</sup> DMM,<sup>61</sup> and FOE<sup>34</sup> formulations of SCF theory. While  $\mathbf{S}^{-1}$  is typically dense, its Cholesky factors  $\mathbf{Z}$  may be quite sparse, depending on the ordering of  $\mathbf{S}$ .<sup>62,63</sup> Thus, if the product  $\mathbf{S}^{-1}\mathbf{X}$  is sparse, it is possible to efficiently apply the inverse as

$$\mathbf{Z}(\mathbf{Z}^T\mathbf{X})=\mathbf{S}^{-1}\mathbf{X}, \quad (12)$$

without ever referencing a potentially dense  $\mathbf{S}^{-1}$ .

### C. The LNVD objective

Beginning with the seminal work of Löwdin<sup>53</sup> and McWeeny,<sup>49</sup> it has long been recognized that eigensolution of the SCF equations can be avoided by the direct minimization of  $E_{\text{el}}$  subject to the constraint of normalization

$$N_{\text{el}}=\text{Tr}\{\mathbf{P}\}, \quad (13)$$

and idempotency,

$$\mathbf{P}=\mathbf{P}\mathbf{P}. \quad (14)$$

This latter condition is equivalent to requiring orthogonality of  $\mathbf{C}_{\text{occ}}$ .

Variation of the electronic energy leads to the stationary condition<sup>49,64</sup>

$$\delta E_{\text{el}}=2\text{Tr}\{\mathbf{F}\delta\mathbf{P}\}=0, \quad (15)$$

allowing  $E_{\text{el}}$  to be replaced with the equivalent objective

$$\mathcal{E}=\text{Tr}\{\mathbf{P}\mathbf{F}\} \quad (16)$$

in which  $\mathbf{F}$  no longer depends on  $\mathbf{P}$ .

A number of workers have formulated different versions of density matrix minimization,<sup>65-71</sup> but these were not competitive with eigensolution due to the expense associated with enforcing Eq. (14), and the high efficiency of direct eigensolvers.<sup>72</sup> A breakthrough came with the density matrix minimization of Li, Nunes, and Vanderbilt<sup>17</sup> and Daw<sup>18</sup> (LNVD), which imposes the constraint on idempotency implicitly through substitution of the McWeeny purification,<sup>49</sup>

$$\mathbf{P}=\mathbf{3}\mathbf{P}\mathbf{P}-\mathbf{2}\mathbf{P}\mathbf{P}\mathbf{P}, \quad (17)$$

into Eq. (16). The purification transform brings an approximately idempotent matrix closer to idempotency, a process that is quadratically convergent upon iteration.<sup>49</sup> With this substitution, Eq. (14) is imposed approximately during optimization, and may be used after convergence to restore idempotency to within the accuracy allowed by the underlying matrix algebra. In the LNVD density matrix minimiza-

tion, normalization is imposed through a Lagrange multiplier  $\mu$ , yielding the LNVD objective

$$\Omega=\text{Tr}\{(\mathbf{3}\mathbf{P}\mathbf{P}-\mathbf{2}\mathbf{P}\mathbf{P}\mathbf{P})\mathbf{F}\}+\mu(N_{\text{el}}-\text{Tr}\mathbf{P}). \quad (18)$$

### III. DENSITY MATRIX MINIMIZATION

In the LNVD method, a density matrix minimization is carried out after forming each new Fock matrix. To achieve an  $\mathcal{O}(N)$  complexity, it is necessary to employ an approximate sparse matrix algebra and to forego methods that require storing a Hessian. One possibility is steepest descent (SD) in which the gradient,

$$\mathbf{G}_j=-\nabla_{\mathbf{P}_j}\Omega \quad (19)$$

is simply followed. A better alternative is the nonlinear conjugate gradient (NLCG) method, which employs properties of the Hessian implicitly through a sequence of  $A$ -orthogonal search directions  $\mathbf{H}_j$ .<sup>73-75</sup> In the nonlinear conjugate gradient method, if an exact line search is performed, that is if a steplength  $\lambda_j$  is chosen such that the updated density matrix

$$\mathbf{P}_{j+1}=\mathbf{P}_j+\lambda\mathbf{H}_j \quad (20)$$

exactly minimizes the objective, then each gradient in the sequence will be orthogonal to the others;<sup>76</sup>

$$(\mathbf{G}_k,\mathbf{G}_j)=0 \quad (k\neq j). \quad (21)$$

If conjugacy is preserved, the nonlinear conjugate gradient method is in general much more efficient than the steepest descent. In practice, incomplete matrix algebra, deviation from quadratic behavior, and inexact line searches can spoil this property,<sup>75-77</sup> leading to an algorithm that may be less efficient than the steepest descent. Note that modifying the gradient or objective during optimization, through a purification step or extrapolation, will tend to destroy conjugacy.

#### A. Gradient and line search

Using the trace algebra,<sup>78</sup> the gradient of  $\Omega$  is

$$\nabla\Omega=3(\mathbf{P}\mathbf{F}+\mathbf{F}\mathbf{P})-2(\mathbf{P}\mathbf{P}\mathbf{F}+\mathbf{P}\mathbf{F}\mathbf{P}+\mathbf{F}\mathbf{P}\mathbf{P})-\mu\mathbf{I}, \quad (22)$$

where it has been assumed that  $\mu$  does not depend on  $\mathbf{P}$ . The Lagrange multiplier  $\mu$  is chosen as

$$\mu=\text{Tr}\{3(\mathbf{P}\mathbf{F}+\mathbf{F}\mathbf{P})-2(\mathbf{P}\mathbf{P}\mathbf{F}+\mathbf{P}\mathbf{F}\mathbf{P}+\mathbf{F}\mathbf{P}\mathbf{P})\}/N_{\text{bas}}, \quad (23)$$

which renders the gradient traceless at each step.<sup>20,23</sup> With each step an analytic line search is carried out as in Refs. 20 and 23 by solving

$$\frac{d\Omega[\mathbf{P}+\lambda\mathbf{H}]}{d\lambda}=b+2c\lambda+3d\lambda^2=0 \quad (24)$$

for the root that minimizes  $\Omega$ , where

$$b=\text{Tr}\{\mathbf{H}\nabla\Omega\} \quad (25)$$

$$c=\text{Tr}\{3\mathbf{H}\mathbf{H}\mathbf{F}-2(\mathbf{H}\mathbf{H}\mathbf{P}\mathbf{F}+\mathbf{H}\mathbf{P}\mathbf{H}\mathbf{F}+\mathbf{P}\mathbf{H}\mathbf{H}\mathbf{F})\} \quad (26)$$

$$d=-2\text{Tr}\{\mathbf{H}\mathbf{H}\mathbf{H}\mathbf{F}\}. \quad (27)$$

## B. A simplifying guess

At the start of each minimization, it is necessary to choose a guess density matrix. A natural choice is the density matrix from the previous SCF cycle. However, this guess leads to very slow convergence, typically requiring on the order of 50–100 nonlinear conjugate gradient cycles.

One method to accelerate convergence is to precondition with a diagonal Hessian,<sup>74,75,79</sup> which has been used by Millam and Scuseria.<sup>23</sup> Another approach employed by Millam and Scuseria is to use direct inversion in the iterative subspace (DIIS)<sup>80</sup> to minimize the commutator  $[\mathbf{F}^i, \mathbf{P}_j^{i+1}]$  between nonlinear conjugate gradient cycles (DIIS-CG), where superscripts have been used to emphasize the SCF iteration, and  $j$  is the CG cycle. Note that  $[\mathbf{F}^j, \mathbf{P}^j]$  will be zero only when  $\mathbf{P}^j = \mathbf{P}^{j+1}$ .

An alternative guess for the start of a density matrix minimization is

$$\mathbf{P}_0 = \left( \frac{N_{\text{el}}}{N_{\text{bas}}} \right) \mathbf{I}, \quad (28)$$

which preserves symmetry of the converged result, namely

$$[\mathbf{P}, \mathbf{F}] = [\mathbf{H}, \mathbf{F}] = 0, \quad (29)$$

throughout a nonlinear conjugate gradient or steepest descent sequence.

Starting with Eq. (28) it is found that after some characteristic number of conjugate gradient cycles,  $N_{\text{cg}} \sim 3$ , purification will establish the exact answer, even though the gradient may still be large. If too few steps are taken, purification will lead to the loss of electrons. As  $N_{\text{cg}}$  increases, fewer purification steps  $N_{\text{pur}}$  are needed to obtain the correct result. In every case, it is found possible to cast the dominant portion of the work into the purification step, which is quadratically convergent.

In addition to enhanced convergence properties, Eq. (29) leads to a significant reduction in complexity of the gradient evaluation and the line search. In particular, Eqs. (22) and (23) reduce to

$$\nabla \Omega = 6(\mathbf{I} - \mathbf{P})\mathbf{P}\mathbf{F} - \mu \mathbf{I} \quad (30)$$

and

$$\mu = 6 \text{Tr}\{(\mathbf{I} - \mathbf{P})\mathbf{P}\mathbf{F}\} / N_{\text{bas}}, \quad (31)$$

while Eq. (26) simplifies to

$$c = 3 \text{Tr}\{(\mathbf{I} - 2\mathbf{P})\mathbf{H}\mathbf{H}\mathbf{F}\}. \quad (32)$$

Inspection of Eq. (30) shows that during minimization,  $\mathbf{P}$  will start with the sparsity pattern of  $\mathbf{F}$ , and will broaden with additional iterations. In practice, this results in initial cycles that are quite-fast.

Note that preconditioning will destroy the commutation property; that is, if preconditioning is used then Eqs. (30), (31), and (32) will become invalid.

## C. Nonorthogonal density matrix minimization

Just as with eigensolution, it is possible to cast the LNVD method in a nonorthogonal basis. There are at least

two ways to perform the density matrix minimization in a nonorthogonal basis: One way is to simply follow the gradient of the nonorthogonal objective<sup>19,21,22</sup>

$$\Omega_{\text{ao}} = \text{Tr}\{(3P_{\text{ao}}\mathbf{S}\mathbf{P}_{\text{ao}} - 2\mathbf{P}_{\text{ao}}\mathbf{S}\mathbf{P}_{\text{ao}}\mathbf{S}\mathbf{P}_{\text{ao}})\mathbf{F}_{\text{ao}}\} - \mu N_{\text{el}}. \quad (33)$$

Another approach is to follow the gradient of the orthogonal objective, Eq. (18), but in a nonorthogonal basis as suggested by White *et al.*,<sup>61</sup>

$$(\nabla \Omega)_{\text{ao}} = \mathbf{S}^{-1/2} [\nabla \Omega] \mathbf{S}^{-1/2}. \quad (34)$$

This latter approach is independent of basis set transformations, and may have improved convergence properties over the former approach. One disadvantage is that the inverse overlap matrix  $\mathbf{S}^{-1}$  must be applied at every iteration, and it is no longer possible to use the nonlinear conjugate gradient without a back transformation, as the gradient corresponds to the orthogonal objective, Eq. (18), rather than the nonorthogonal one, Eq. (33). On the other hand,  $\mathbf{F}_{\text{ao}}$  tends to be much sparser than  $\mathbf{F}$ , and it may be that working in the AO results in faster gradient evaluations and line minimizations.

Starting with the guess,

$$\mathbf{P}_{\text{ao},0} = \left( \frac{N_{\text{el}}}{N_{\text{bas}}} \right) \mathbf{S}^{-1} \quad (35)$$

it is possible to achieve simplification and improved convergence properties in the AO basis that are equivalent to those of the orthogonal simplified density matrix minimization developed in Sec. IIIB. With this guess, Eq. (34) reduces to

$$(\nabla \Omega)_{\text{ao}} = 6\mathbf{S}^{-1} [(\mathbf{I} - \mathbf{S}\mathbf{P}_{\text{ao}})\mathbf{F}_{\text{ao}}\mathbf{P}_{\text{ao}} - \mu \mathbf{I}]. \quad (36)$$

with

$$\mu = \text{Tr}\{(\mathbf{I} - \mathbf{S}\mathbf{P}_{\text{ao}})\mathbf{F}_{\text{ao}}\mathbf{P}_{\text{ao}}\} / N_{\text{bas}}. \quad (37)$$

Rather than compute and manipulate  $\mathbf{S}^{-1}$ , it is much more efficient to apply the AINV factors as in Eq. (12). As before, an analytic line search is carried out. However, this time the nonorthogonal objective is used, yielding the coefficients

$$b = 6 \text{Tr}\{(\mathbf{I} - \mathbf{P}_{\text{ao}}\mathbf{S})\mathbf{H}\mathbf{F}_{\text{ao}}\}, \quad (38)$$

$$c = 3 \text{Tr}\{(\mathbf{I} - 2\mathbf{P}_{\text{ao}}\mathbf{S})\mathbf{H}\mathbf{S}\mathbf{H}\mathbf{F}_{\text{ao}}\}, \quad (39)$$

$$d = -2 \text{Tr}\{\mathbf{H}\mathbf{S}\mathbf{H}\mathbf{S}\mathbf{H}\mathbf{F}_{\text{ao}}\}. \quad (40)$$

## IV. IMPLEMENTATION

The orthogonal and nonorthogonal simplified density matrix minimization have been implemented in the MondoSCF<sup>52</sup> suite of programs for linear scaling SCF theory, with details as follows.

### A. Atom-blocked sparse matrix algebra

A library of routines for carrying out variable block size, sparse matrix algebra have been developed and employed in the implementation of MondoSCF. These routines exploit both existing sparsity and data locality. Related methods have been shown to yield significant speedups for matrix-vector multiplies.<sup>81–83</sup> Note that in Ref. 82, Goedecker *et al.* employ an atom-blocked algebra for matrix-vector multiplies in the context of the tight binding FOE method.

## 1. Matrix multiplies

In the sparse atom-blocked matrix multiply, the matrices are stored in a modified compressed sparse row (CSR) format,<sup>84–86</sup> in which all pointers index blocks rather than individual matrix elements. An additional pointer is required to address the start of each block since the blocks are not grouped by dimension. A modified version of Gustavson's sparse matrix multiply<sup>87</sup> is employed, with DGEMM-like routines to carry out the block–block multiplies. DGEMM is a level three BLAS<sup>72</sup> routine for performing double precision matrix multiplies. In the limit that the block size is reduced to  $1 \times 1$ , an algorithm equivalent to the standard<sup>87</sup> sparse matrix multiply is obtained.

## 2. Truncation

Small matrix elements are dropped after each matrix operation such as multiplication or addition. Blocks are dropped when they meet the criterion

$$\|C_{IK}\|_F < \text{TrixtNeglect}, \quad (41)$$

where  $\text{TrixtNeglect}$  is a drop tolerance and  $F$  demarks the Frobenius norm<sup>88</sup>

$$\|C_{IK}\|_F = \sqrt{\text{Tr}\{C_{IK}^T C_{IK}\}}. \quad (42)$$

This approach should be compared with truncation of individual elements,<sup>23,89</sup> and separation-based truncation.<sup>17,22</sup> With truncation at the level of elements, a change of basis may lead to very different behavior of errors. This is because the average magnitude of the matrix elements must decrease with increasing basis set size.<sup>40</sup> Truncation at the level of atoms is expected to yield errors that are more transferable, as the amount of charge in a block should be less sensitive to a change of basis.

In separation-based truncation, elements of a matrix product are computed only if the corresponding basis function separation is less than a cut-off radius  $R_c$ . This has the advantage that the structure of all matrices is known *a priori*. On the other hand, this approach requires a careful system, basis set, and model chemistry dependent parametrization. Also, when examining the fall off of density matrix elements with basis function separation,<sup>14</sup> it is seen that many small elements correspond to small separations. This suggests that, for a given level of error, truncation schemes based on magnitude may yield sparser matrices.

## 3. SPAMM

A method related to thresholding schemes used in the order  $N$  exchange method ONX<sup>8,14</sup> has been extended to matrix multiplication, and may be well suited for matrices with elements that decay away from the diagonal. In this sparse approximate matrix multiply (SPAMM), small matrix elements are treated in a more approximate manner than are the large ones. For the matrix multiply  $C = AB$ , each block  $C_{IK}$  is formed as the sum over blocks

$$C_{IK} = \sum_{J=1}^M A_{IJ} B_{JK}. \quad (43)$$

The error

$$\|\delta C_{IK}\|_F < \text{MultNeglect} \quad (44)$$

may be rigorously bounded while avoiding all contractions  $A_{IJ} B_{JK}$  for which

$$\|B_{JK}\|_F < \frac{\text{MultNeglect}}{M \|A_{IJ}\|_F}, \quad (45)$$

where  $M$  is the number of blocks in the contraction, and  $\text{MultNeglect}$  is a threshold. Just as in the case of ONX, repeatedly executing a conditional like Eq. (45) in an innermost loop can create a significant overhead. This is avoided in SPAMM by sorting each row in decreasing order on the basis of  $\|\cdot\|_F$ , and using a binary search to find the limit ( $\leq M$ ) of the innermost loop.

## 4. Trace

A substantial savings may be achieved in the execution of  $\text{Tr}\{\cdot\}$  by performing both the contraction and the diagonal accumulation at the same time as

$$\text{Tr}\{AB\} = \sum_{ik} A_{ik} B_{ki}, \quad (46)$$

rather than first forming the product  $AB$ .

## 5. AINV

MondoSCF employs a blocked version of AINV, which is described in detail elsewhere.<sup>60</sup>

## B. SCF

### 1. Initial guess

For a minimal basis such as STO-3G, a sparse linear scaling guess is built from the superposition of spherically averaged atomic density matrices in an orthogonal (diagonal) representation. The AINV factor  $Z$  is then used to transform to an AO representation.<sup>90</sup> For more complicated basis sets, a partially converged density matrix from a minimal basis set calculation is used to construct a mixed basis Fock matrix.<sup>90</sup> These guesses are linear scaling.

### 2. Eigensolver

The Roothaan equations, Eq. (6), are solved with the LAPACK divide and conquer eigensolver DSYEVD<sup>72</sup> as implemented in the SGI scientific library SCSL. This a competitive algorithm that has been highly optimized for the SGI platform.

### 3. Fock builds

The multipole accelerated symmetrized order  $N$  exchange algorithm (MASONX)<sup>16</sup> is used to compute the exchange matrix, and the quantum chemical tree code (QCTC)<sup>13</sup> is used to compute the Coulomb matrix. Tight integral ( $\text{TwoENeglect} = 10^{-9}$ ) and distribution ( $\text{DistNeglect} = 10^{-11}$ ) thresholds are used throughout.

### 4. DIIS

A sparse implementation of DIIS for extrapolation of the Fock matrix<sup>91</sup> is used in all calculations. In this implementation, the error vector

$$e_{ao} = F_{ao} P_{ao} S - SP_{ao} F_{ao} \quad (47)$$

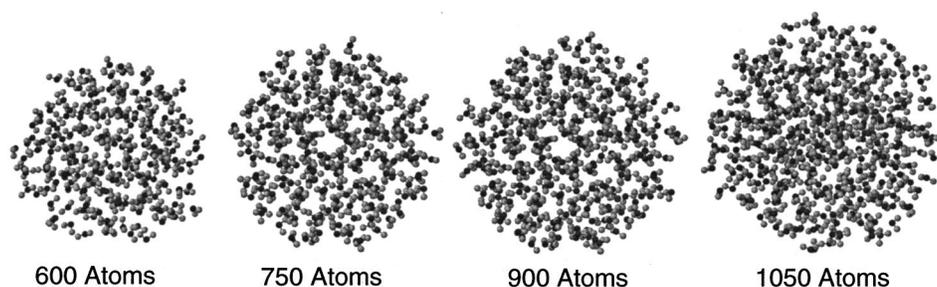


FIG. 1. Configuration of the clusters with 200, 250, 300 and 350 water molecules.

is transformed into an orthonormal basis with the transformation

$$\mathbf{e} = \mathbf{Z}^T \mathbf{e}_{\text{ao}} \mathbf{Z}. \quad (48)$$

The use of sparse matrix algebra will eventually cause the DIIS method to stagnate. Tightening up the matrix thresholds with SCF convergence leads to premature stagnation that is incommensurate with the thresholds. One solution is to simply restart the DIIS after each decrement of a threshold, and this should lead to significantly reduced average CPU times per SCF cycle. For clarity and simplicity, this approach has not been followed here; rather, constant thresholds are used throughout.

### 5. Nonlinear conjugate gradient

The Polak–Ribière version of the nonlinear conjugate gradient<sup>92</sup> has been implemented for the orthogonal minimization. As discussed in Sec. III, the nonquadratic behavior of the objective, approximate arithmetic, and deviation of the gradient from the objective due to the normalization constraint may destroy conjugacy as measured by

$$\alpha = \frac{|(\mathbf{G}_j, \mathbf{G}_{j+1})|}{(\mathbf{G}_{j+1}, \mathbf{G}_{j+1})}, \quad (49)$$

which is ideally 0. While the Polak–Ribière nonlinear conjugate gradient tends to restart along  $\mathbf{G}_{j+1}$  when conjugacy is lost, Powell<sup>93</sup> suggests a complete steepest descent restart when  $\alpha$  exceeds 0.2. In the calculations performed here, this condition typically occurs only for large values of  $N_{\text{cg}}$ , and no advantage has been found in performing a steepest descent restart. However, for systems with a small HOMO–LUMO gap this may yield an advantage.

### 6. Overview of the orthogonal method

After each Fock build, DIIS extrapolation of the Fock matrix is used, and the extrapolated Fock matrix is transformed to an orthogonal basis. A diagonal guess is used to start, and  $N_{\text{cg}}$  nonlinear conjugate gradient steps are taken using the simplified gradient and line minimization particular to the SDMM. Then, the resulting density matrix is brought to idempotency through purification. The degree of idempotency achieved depends on the matrix thresholds *TrixNeglect* and *MultNeglect*. Stagnation manifests itself in an inability to further reduce the maximum block of the difference density  $P_{I_{\text{pur}}} - P_{I_{\text{pur}+1}}$ , and is said to occur after  $N_{\text{pur}}$  iterations. Thus,  $N_{\text{pur}}$  is a function of  $N_{\text{cg}}$  and the matrix thresholds. The resulting density matrix is then transformed back to a nonorthogonal atomic orbital representation and used to construct a new Fock matrix. Stagnation of the DIIS procedure also occurs with more approximate matrix thresholds, and the SCF procedure is halted after the DIIS error fails to improve.

## V. RESULTS AND DISCUSSION

All computations were carried out on ASCI Bluemountain, which is a large collection of 195 MHz SGI Origin 2000s. The calculations were all performed in serial and in a shared memory environment (nondedicated).

### A. Benchmark systems

Two three-dimensional benchmark suites are used in this study: a sequence of water clusters and a sequence of estane polymers. Linear systems of polyglycine chains, alkanes, and

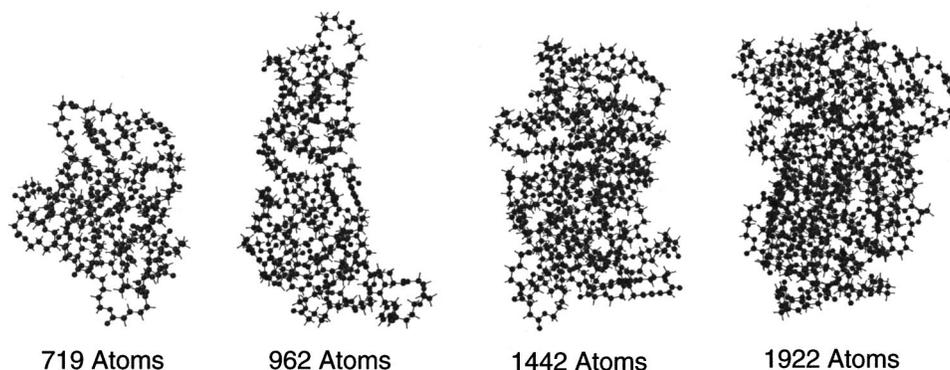


FIG. 2. Local minimum energy configurations of the estane polymers with three, four, six and eight segments, corresponding to 719, 962, 1442, and 1922 atoms, respectively.

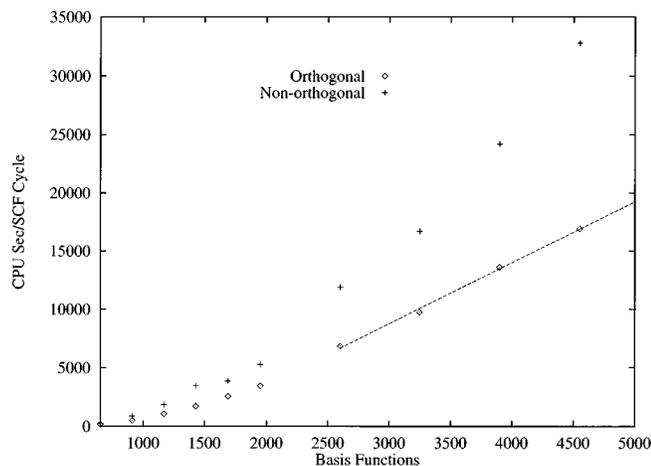


FIG. 3. Scaling of CPU time per SCF for the RHF/3-21G sequence of water clusters using the nonorthogonal and orthogonal simplified density matrix minimization.

nucleic acids yield overly optimistic results with respect to both computational complexity and error control, and are not considered here.

The sequence of clusters includes 50, 70, 90, 110, 150, 200, 250, 300, and 350 water molecules, with a constant density corresponding to standard temperature and pressure. The four largest clusters are shown in Fig. 1. These water clusters have been used in a number of other studies,<sup>7,6,13,10,14,23,89,15,16</sup> and have proven a challenge to achieving linear scaling in computation of the Coulomb matrix,<sup>7,6,10</sup> the exchange matrix,<sup>10</sup> and the density matrix.<sup>23,89</sup>

Estane 5703 polyester urethane is a segmented copolymer which is a constituent in the plastic-bonded high explosive PBX-9501.<sup>94</sup> A fundamental repeat unit,  $N_{\text{seg}}=1$ , is a length of chain comprised of two 4,4'-diphenylmethane diisocyanate and 1,4-butanediol segments bonded to five poly tetramethylene adipate segments. Configurations for degrees of polymerization corresponding to 1, 2, 3, 4, 6, and 8 (242, 482, 719, 962, 1442, and 1922 atoms, respectively) have been generated by molecular dynamics equilibration followed by energy minimization. The four largest configurations are shown in Fig. 2.

## B. Scaling

### 1. Orthogonal and nonorthogonal simplified density matrix minimization

In Fig. 3, the scaling of the nonorthogonal and orthogonal simplified density matrix minimization are shown for the sequence of water clusters at the RHF/3-21G level of theory. In both cases, the thresholds  $\text{TrixNeglect}=10^{-6}$  and  $\text{MultNeglect}=10^{-8}$  were used. Seven minimization steps were taken in both the steepest descent nonorthogonal and the nonlinear conjugate gradient orthogonal case. The non-orthogonal simplified density matrix minimization is much slower than the orthogonal one, and does not appear to achieve linear scaling.

While  $\mathbf{F}$  is denser than  $\mathbf{F}_{\text{ao}}$ , the interspersal of  $\mathbf{S}$  in the nonorthogonal expressions leads to intermediate matrices,

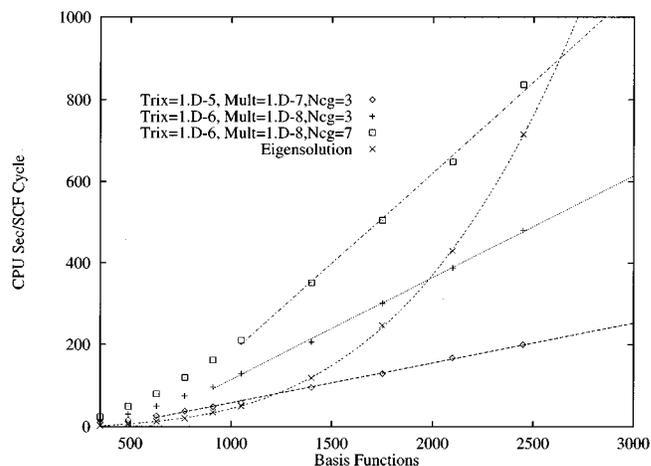


FIG. 4. Scaling of CPU time per SCF for the RHF/STO-3G sequence of water clusters.

e.g.,  $\mathbf{S}\mathbf{F}_{\text{ao}}$ , with an increased density. In addition, the nonorthogonal simplified density matrix minimization requires many more matrix multiplies than does the orthogonal version.

## 2. Water

The scaling of the orthogonal simplified density matrix minimization is shown in Fig. 4 for the water clusters at the RHF/STO-3G level of theory. A factor of two reduction in CPU time is observed on going from  $N_{\text{cg}}=7$  to  $N_{\text{cg}}=3$  for the  $\text{TrixNeglect}=10^{-6}$  calculations. A similar gain is found (but not shown) for the  $\text{TrixNeglect}=10^{-5}$  calculation.

These results should be compared with Fig. 5 from Ref. 23, which shows the inability of the CG-DMS algorithm to achieve linear scaling for LDA/STO-3G calculations on the same set of water clusters. In that work, the onset of linear scaling is predicted to occur at about 400 water molecules. Here, the onset of linear scaling is observed at 90 water molecules for the  $\text{TrixNeglect}=10^{-5}$  calculation, and at 130 water molecules for the  $\text{TrixNeglect}=10^{-6}$  calculations. In addition to differences in the algorithms and implementation, it is known that the LDA tends to underestimate the HOMO-LUMO gap, while HF tends to overestimate it. This effect is shown in Table I. As the asymptotic behavior of the density

TABLE I. HOMO and LUMO energies and the HOMO-LUMO gap for a cluster of 30 waters. All values were obtained with GAUSSIAN 94.<sup>a</sup>

| Theory | Basis set | HOMO   | LUMO   | $E_{\text{gap}}$ |
|--------|-----------|--------|--------|------------------|
| HF     | STO-3G    | -0.318 | 0.497  | 0.815            |
| HF     | 3-21G     | -0.393 | 0.164  | 0.557            |
| HF     | 6-31G**   | -0.423 | 0.126  | 0.549            |
| LDA    | STO-3G    | 0.005  | 0.200  | 0.195            |
| LDA    | 3-21G     | -0.133 | -0.039 | 0.094            |
| LDA    | 6-31G**   | -0.190 | -0.061 | 0.129            |
| B3LYP  | 3-21G     | -0.178 | 0.008  | 0.186            |
| B3LYP  | 6-31G**   | -0.226 | -0.013 | 0.213            |

<sup>a</sup>See Ref. 95.

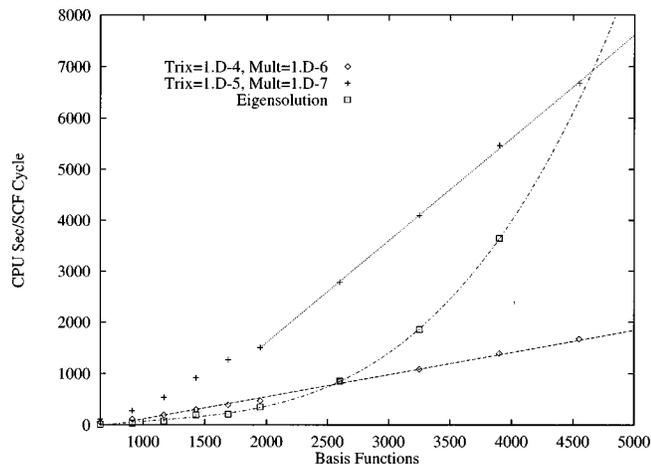


FIG. 5. Scaling of CPU time per SCF for the RHF/3-21G sequence of water clusters.

matrix is governed by  $\sqrt{E_{\text{gap}}}$  as in Eq. (1), this phenomenon may explain some of the differences between this work and Ref. 23.

The scaling of the orthogonal simplified density matrix minimization with  $N_{\text{cg}}=3$  is shown in Figs. 5 and 6 for the water series at the RHF/3-21G level of theory. For the 3-21G calculations, the differences in CPU time between  $N_{\text{cg}}=3$  and  $N_{\text{cg}}=7$  are much less pronounced than in the case of the STO-3G calculations. This is shown in Fig. 6 for  $\text{TrixNeglect}=10^{-6}$ . A similar behavior is also found for looser values.

The onset of linear scaling is slower for the 3-21G calculations than for the STO-3G calculations, occurring at 90, 150, and 200 water molecules for  $\text{TrixNeglect}=10^{-4}$ ,  $10^{-5}$ , and  $10^{-6}$  respectively. Also, the STO-3G density matrix is less dense than the 3-21G density matrix. These differences are most likely because the STO-3G band gap is larger than the 3-21G band gap, as shown in Table I.

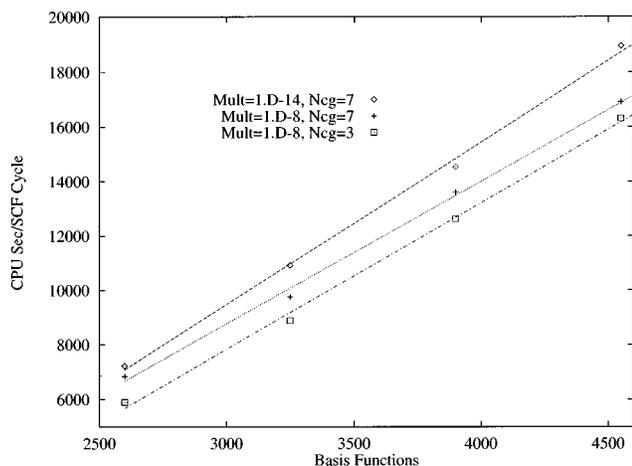


FIG. 6. Scaling of CPU time per SCF for the RHF/3-21G sequence of water clusters with  $\text{TrixNeglect}=10^{-6}$ .

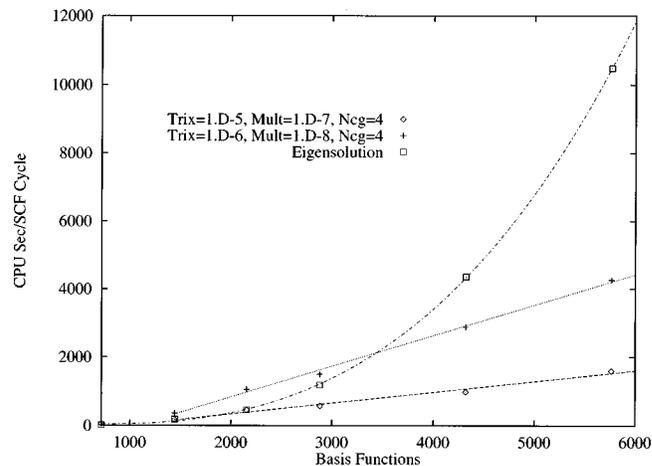


FIG. 7. Scaling of CPU time per SCF for the RHF/STO-3G sequence of estane polymers.

### 3. Estane

The scaling of the orthogonal simplified density matrix minimization with  $N_{\text{cg}}=4$  is shown in Fig. 7 for the estane polymers at the RHF/STO-3G level of theory. As in the water calculations, the  $N_{\text{cg}}=4$  results are about a factor of two faster than those obtained with  $N_{\text{cg}}=8$  (not shown). In all cases, the onset of linear scaling is seen to occur at the  $N_{\text{seg}}=2$  polymer, which corresponds to 482 atoms and 722 basis functions.

### C. Errors

Errors in converged total energies are shown in Fig. 8 for the RHF/STO-3G water series, in Fig. 9 for the RHF/3-21G water series, and in Fig. 10 for the RHF/STO-3G estane sequence. A striking feature of the water errors is that the calculations with  $N_{\text{cg}}=3$  yield more accurate results than the calculations with  $N_{\text{cg}}=7$ .

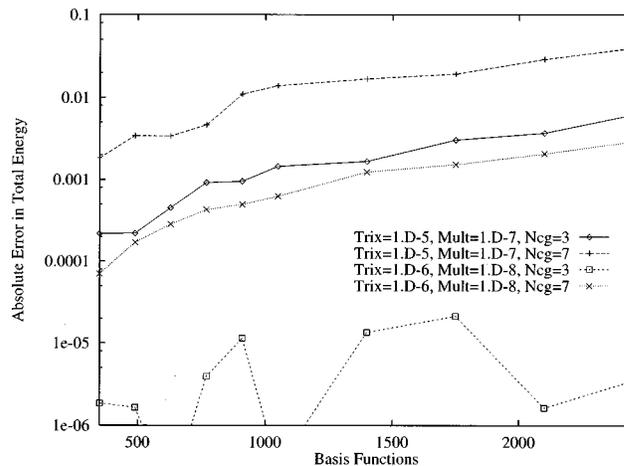


FIG. 8. Absolute errors in converged total energies for the RHF/STO-3G sequence of water clusters.

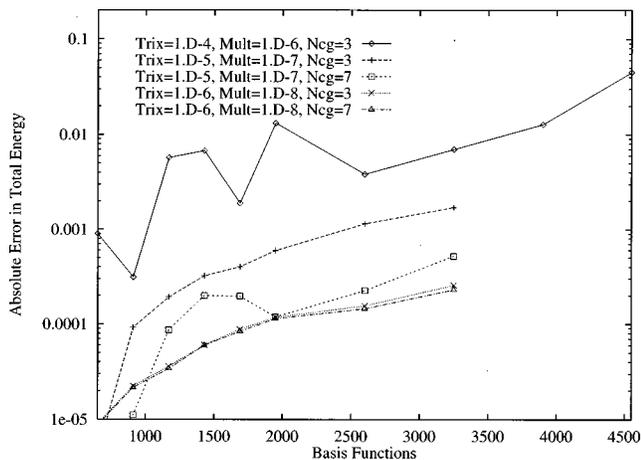


FIG. 9. Absolute errors in converged total energies for the RHF/3-21G sequence of water clusters.

One explanation for this behavior is that, with approximate algebra,  $[F,P]$  and  $[F,H]$  are not exactly zero. With each additional step these commutators grow, leading to errors in the gradient and line search. As the matrix thresholds are tightened, results for different values of  $N_{cg}$  will become identical; this can be seen for the  $\text{TrixNeglect}=10^{-6}$  calculations in Fig. 9.

## D. Performance of the matrix multiply

### 1. MFLOPS

Megafloating point operations per second (MFLOPS) for gradient evaluation (GradE), line minimization (LineM), and purification (Purify) are given in Table II for water, estane, and fullerene for the STO-3G and 3-21G basis sets. These are sustained values, and while dominated by multiplication, they include additions, trace, diagonal adds, matrix-scalar multiplication and input/output (IO). For reference, Warner's<sup>96</sup> optimized DGEMM using  $50 \times 50$  subblocks achieves a peak rate of approximately 270 MFLOPS on the 195 MHz SGI Origin 2000.

The MFLOPS rate for line minimization is less than for purification or gradient evaluation. This is because cache and

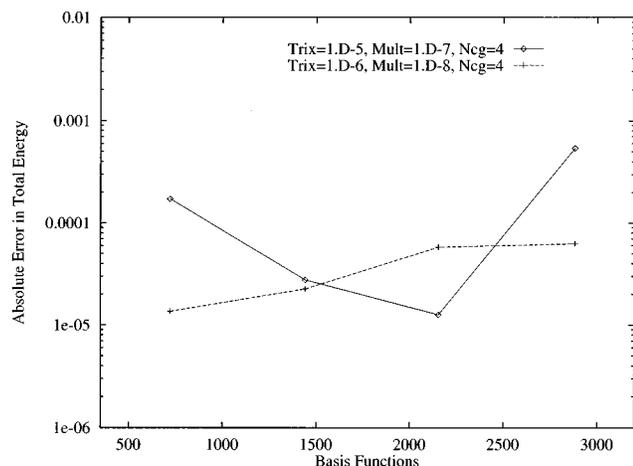


FIG. 10. Absolute errors in converged total energies for the RHF/STO-3G sequence of estane polymers.

TABLE II. Sustained MFLOPS rates as a function of different systems and basis sets.  $N_H/N_X$  is the ratio of hydrogen atoms to heavy atoms.

| System    | $N_H/N_X$ | Basis set | MFLOPS |       |        |
|-----------|-----------|-----------|--------|-------|--------|
|           |           |           | GradE  | LineM | Purify |
| Water     | 2         | STO-3G    | 71     | 57    | 81     |
| Estane    | 1         | STO-3G    | 100    | 83    | 104    |
| Fullerene | 0         | STO-3G    | 162    | 136   | 158    |
| Water     | 2         | 3-21G     | 137    | 125   | 140    |
| Estane    | 1         | 3-21G     | 154    | 135   | 150    |
| Fullerene | 0         | 3-21G     | 163    | 145   | 157    |

integer overheads are higher per FLOP for the evaluation of  $\text{Tr}\{\cdot\}$  with Eq. (46) than for the matrix multiply and addition.

As the ratio of hydrogens to heavy atoms  $N_H/N_X$  decreases, the MFLOPS rate increases. This is because the number of basis functions per atom is smaller (one for STO-NG and two for 3-21G) for hydrogen than for the heavy atoms (five for STO-NG and nine for 3-21G), and computations involving small dimensions are inefficient. This effect is less pronounced for the 3-21G basis set because the CPU time in this case is dominated by the  $9 \times 9 \times 9$  DGEMM corresponding to the all heavy atom multiplies. As larger basis sets are used, the influence of  $N_H/N_X$  will become less significant, and the peak rate will be determined by DGEMM efficiency, which is limited by cache effects and is block size dependent.<sup>96,97</sup>

### 2. SPAMM

The choice

$$\text{MultNeglect} = 10^{-2} \text{TrixNeglect} \quad (50)$$

has been found to yield results that are very close to those obtained with  $\text{MultNeglect}=0$ , and this parametrization has been used throughout. In Fig. 6, results are shown using SPAMM with this default parametrization and with  $\text{MultNeglect}=10^{-14}$ . The effects of SPAMM here and in the other systems studied are small ( $<10\%$ ). Nevertheless, it may be more effective in systems which are not as well localized as those studied here.

### 3. Nonlinear conjugate gradient vs purification cycles

As mentioned in Sec. III B, the CPU time per nonlinear conjugate gradient step grows with the number of steps taken. This is shown in Table III for the 3-21G 350 water

TABLE III. The increase in CPU time per orthogonal nonlinear conjugate gradient cycle for the 3-21G 350 water calculation with  $\text{TrixNeglect}=10^{-7}$  and  $\text{MultNeglect}=10^{-5}$ .

| $I_{CG}$ | CPU (s) |
|----------|---------|
| 0        | 407     |
| 1        | 505     |
| 2        | 714     |
| 3        | 780     |
| 4        | 749     |
| 5        | 905     |
| 6        | 828     |

TABLE IV. Decrease in the number of purification cycles  $N_{\text{pur}}$  on going from  $N_{\text{cg}}=3$  to  $N_{\text{cg}}=7$  for the 3-21G 350 water calculation with  $\text{TrixNeglect}=10^{-5}$  and  $\text{MultNeglect}=10^{-7}$ .

| $I_{\text{pur}}$ | $\text{Max}\ \Delta\mathbf{P}_{IJ}\ _F$ |                   |
|------------------|---|-------------------|
|                  | $N_{\text{cg}}=3$                       | $N_{\text{cg}}=7$ |
| 1                | 0.215D+00                               | 0.147D+0          |
| 2                | 0.154D+00                               | 0.760D-01         |
| 3                | 0.979D-01                               | 0.434D-01         |
| 4                | 0.662D-01                               | 0.227D-01         |
| 5                | 0.579D-01                               | 0.702D-02         |
| 6                | 0.411D-01                               | 0.793D-03         |
| 7                | 0.282D-01                               | 0.327D-04         |
| 8                | 0.175D-01                               | 0.328D-04         |
| 9                | 0.813D-02                               | ...               |
| 10               | 0.182D-02                               | ...               |
| 11               | 0.954D-04                               | ...               |
| 12               | 0.341D-04                               | ...               |
| 13               | 0.330D-04                               | ...               |

calculation. Note that the first cycle's cost is half as much as the last cycle. The reduction in the number of purification steps required to reach stagnation is shown in Table IV for  $N_{\text{cg}}=3$  and  $N_{\text{cg}}=7$ . Stagnation occurs when  $\text{Max}\|\Delta\mathbf{P}_{IJ}\|_F$  fails to decrease with continued iteration due to the approximate matrix algebra.

## VI. CONCLUSIONS

A simplified density matrix minimization (SDMM) has been introduced and shown to yield linear scaling CPU times that are highly competitive with a state-of-the-art dense eigensolver. This work is the first demonstration of linear scaling HF theory for three-dimensional systems, and one of only a few that consider double zeta basis sets or the behavior of errors with increasing system size.

The orthogonal simplified density matrix minimization uses four fewer matrix multiplies than conventional implementations of the LNVD density matrix minimization, and requires no transpositions. In addition, it is found that only a few minimization steps are required, and that the majority of the work can be shifted to the quadratically convergent McWeeny purification. For more approximate matrix thresholds, taking fewer nonlinear conjugate gradient steps is found to yield both faster CPU times and more accurate results.

The nonorthogonal simplified density matrix minimization turns out to be noncompetitive with the orthogonal version, and is unable to achieve linear scaling.

The AINV algorithm for computation of the inverse Cholesky factor  $\mathbf{Z}$  has been introduced to linear scaling electronic structure theory, and found to be essential for working with nonorthogonal basis sets. In particular,  $\mathbf{Z}$  is used in the transformation to and from an orthogonal representation, formation of a linear scaling initial guess, application of the inverse overlap, and in DIIS extrapolation.

A sparse atom-blocked matrix algebra was introduced and found to yield sustained performances of up to 160 MFLOPS, which is more than half of the 270 MFLOPS possible with an optimized DGEMM. It is interesting to con-

sider an extension of this approach to uniform block sizes that are cache optimal, which could yield a factor of two speedup.

The methods developed here should lead to very efficient parallel implementations. Sparse matrix multiplication is efficient in parallel, and blocking will improve this efficiency by reducing the overhead associated with decomposition and reordering. Also, because the computation of  $\mathbf{Z}$  reduces to performing matrix multiplies rather than triangular solves (as in the case of  $\mathbf{L}$ ), the AINV algorithm is ideally suited for parallel implementations.

Linear scaling Hartree-Fock theory for insulating three-dimensional systems has been clearly established in Refs. 13 and 14 for the Fock build, and here for solution of the SCF equations. While achieving linear scaling may be more difficult for model chemistries that yield a less inflated HOMO-LUMO gap, with parallel implementation and code maturation the outlook for large scale application of the simplified density matrix minimization is bright.

## ACKNOWLEDGMENTS

The ASCI project at LANL is gratefully acknowledged for support of this work. The Advanced Computing Laboratory of Los Alamos National Laboratory, Los Alamos, NM is acknowledged. This work was performed on computing resources located at this facility. Thanks to Michele Benzi for help with AINV, and for many enlightening conversations. Thomas D. Sewell is thanked for furnishing the estane coordinates. Thanks also to Yousef Saad for suggesting the use of a sparse approximate inverse in nonorthogonal minimization.

- <sup>1</sup>D. L. Strout and G. E. Scuseria, *J. Chem. Phys.* **102**, 8448 (1995).
- <sup>2</sup>J. Almlöf, K. Faegri, and K. Korsell, *J. Comput. Chem.* **3**, 385 (1982).
- <sup>3</sup>M. Häser and R. Ahlrichs, *J. Comput. Chem.* **10**, 104 (1989).
- <sup>4</sup>C. A. White, B. Johnson, P. Gill, and M. Head-Gordon, *Chem. Phys. Lett.* **230**, 8 (1994).
- <sup>5</sup>C. A. White, B. G. Johnson, P. M. W. Gill, and M. Head-Gordon, *Chem. Phys. Lett.* **253**, 268 (1996).
- <sup>6</sup>M. Challacombe, E. Schwegler, and J. Almlöf, *J. Chem. Phys.* **104**, 4685 (1996).
- <sup>7</sup>M. Challacombe, E. Schwegler, and J. Almlöf, in *Computational Chemistry: Review of Current Trends*, edited by J. Leszczynski (World Scientific, Singapore, 1996), pp. 53–107.
- <sup>8</sup>E. Schwegler and M. Challacombe, *J. Chem. Phys.* **105**, 2726 (1996).
- <sup>9</sup>J. P. Dombroski, S. W. Taylor, and P. M. W. Gill, *J. Phys. Chem.* **100**, 6272 (1996).
- <sup>10</sup>J. C. Burant, R. E. Stratmann, and M. J. Frisch, *J. Chem. Phys.* **105**, 8969 (1996).
- <sup>11</sup>R. E. Stratmann, G. E. Scuseria, and M. J. Frisch, *Chem. Phys. Lett.* **257**, 213 (1996).
- <sup>12</sup>M. C. Strain, G. E. Scuseria, and M. J. Frisch, *Science* **271**, 51 (1996).
- <sup>13</sup>M. Challacombe and E. Schwegler, *J. Chem. Phys.* **106**, 5526 (1997).
- <sup>14</sup>E. Schwegler, M. Challacombe, and M. Head-Gordon, *J. Chem. Phys.* **106**, 9708 (1997).
- <sup>15</sup>C. Ochsenfeld, C. A. White, and M. Head-Gordon, *J. Chem. Phys.* **109**, 1663 (1998).
- <sup>16</sup>E. Schwegler and M. Challacombe (unpublished).
- <sup>17</sup>X. P. Li, R. W. Nunes, and D. Vanderbilt, *Phys. Rev. B* **47**, 10891 (1993).
- <sup>18</sup>M. S. Daw, *Phys. Rev. B* **47**, 10895 (1993).
- <sup>19</sup>R. W. Nunes and D. Vanderbilt, *Phys. Rev. B* **50**, 17611 (1994).
- <sup>20</sup>S. Y. Qiu, C. Z. Wang, K. M. Ho, and C. T. Chan, *J. Phys.: Condens. Matter* **6**, 9153 (1994).
- <sup>21</sup>E. Hernández and M. J. Gillan, *Phys. Rev. B* **51**, 10157 (1995).
- <sup>22</sup>E. Hernández, M. J. Gillan, and C. Goringe, *Phys. Rev. B* **53**, 7147 (1996).
- <sup>23</sup>J. M. Millam and G. E. Scuseria, *J. Chem. Phys.* **106**, 5569 (1997).

- <sup>24</sup>C. Ochsenfeld and M. Head-Gordon, *Chem. Phys. Lett.* **270**, 399 (1997).
- <sup>25</sup>F. Mauri, G. Galli, and R. Car, *Phys. Rev. B* **47**, 9973 (1993).
- <sup>26</sup>P. Ordejón, D. A. Drabold, M. P. Grumbach, and R. M. Martin, *Phys. Rev. B* **48**, 14646 (1993).
- <sup>27</sup>J. Kim, F. Mauri, and G. Galli, *Phys. Rev. B* **51**, 1456 (1995).
- <sup>28</sup>D. Sanchezportal, P. Ordejón, E. Artacho, and J. M. Soler, *Int. J. Quantum Chem.* **65**, 453 (1997).
- <sup>29</sup>P. Ordejón, E. Artacho, and J. M. Soler, *Phys. Rev. B* **53**, 10441 (1996).
- <sup>30</sup>S. Goedecker and L. Colombo, *Phys. Rev. Lett.* **73**, 122 (1994).
- <sup>31</sup>S. Goedecker, *J. Comput. Chem.* **118**, 261 (1995).
- <sup>32</sup>Y. Huang, D. J. Kouri, and D. K. Hoffman, *Chem. Phys. Lett.* **243**, 367 (1995).
- <sup>33</sup>A. F. Voter, J. D. Kress, and R. N. Silver, *Phys. Rev. B* **53**, 12733 (1996).
- <sup>34</sup>H. Roder, R. N. Silver, D. A. Drabold, and J. J. Dong, *Phys. Rev. B* **55**, 15382 (1997).
- <sup>35</sup>F. Gagel, *J. Comput. Chem.* **139**, 399 (1998).
- <sup>36</sup>G. Galli, *Curr. Opin. Solid State Mater. Sci.* **1**, 864 (1996).
- <sup>37</sup>D. R. Bowler *et al.*, *Modell. Simul. Mater. Sci. Eng.* **5**, 199 (1997).
- <sup>38</sup>S. Goedecker, *Rev. Mod. Phys.* (submitted).
- <sup>39</sup>W. Kohn, *Int. J. Quantum Chem.* **56**, 229 (1995).
- <sup>40</sup>P. Maslen *et al.*, *J. Phys. Chem. A* **102**, 2215 (1998).
- <sup>41</sup>E. Hernández, M. J. Gillan, and C. M. Goringe, *Phys. Rev. B* **55**, 13485 (1997).
- <sup>42</sup>M. Lepetit, L. Lafon, and X. Lafage, *Int. J. Quantum Chem.* **64**, 411 (1997).
- <sup>43</sup>S. Q. Wei and M. Y. Chou, *Phys. Rev. Lett.* **76**, 2650 (1996).
- <sup>44</sup>C. J. Tymczak and X. Q. Wang, *Phys. Rev. Lett.* **78**, 3654 (1997).
- <sup>45</sup>R. A. Lippert, T. A. Arias, and A. Edelman, *J. Comput. Chem.* **140**, 278 (1998).
- <sup>46</sup>J. R. Chelikowsky, N. Troullier, and Y. Saad, *Phys. Rev. Lett.* **72**, 1240 (1994).
- <sup>47</sup>J. R. Chelikowsky, N. Troullier, K. Wu, and Y. Saad, *Phys. Rev. B* **50**, 11355 (1994).
- <sup>48</sup>E. L. Briggs, D. J. Sullivan, and J. Bernholc, *Phys. Rev. B* **52**, R5471 (1995).
- <sup>49</sup>R. McWeeny, *Rev. Mod. Phys.* **126**, 1028 (1962).
- <sup>50</sup>M. Benzi and C. D. Meyer, *SIAM J. Sci. Comput.* **16**, 1159 (1995).
- <sup>51</sup>M. Benzi, C. D. Meyer, and M. Tüma, *SIAM J. Sci. Comput.* **17**, 1135 (1996).
- <sup>52</sup>M. Challacombe and E. Schwegler, *mondoSCF a suite of programs for linear scaling SCF theory* (unpublished).
- <sup>53</sup>P. O. Löwdin, *Phys. Rev.* **97**, 1490 (1955).
- <sup>54</sup>A. Szabo and N. S. Ostlund, *Modern Quantum Chemistry*, 1st revised ed. (Mc Graw-Hill, New York, 1989).
- <sup>55</sup>C. C. J. Roothaan, *Rev. Mod. Phys.* **23**, 69 (1951).
- <sup>56</sup>J. H. Wilkinson, *The Algebraic Eigenvalue Problem* (Clarendon, Oxford, 1965).
- <sup>57</sup>G. W. Stewart, *Introduction to Matrix Computations* (Academic, London, 1973).
- <sup>58</sup>P. O. Löwdin, *J. Chem. Phys.* **18**, 365 (1950).
- <sup>59</sup>P. O. Löwdin, *Adv. Phys.* **5**, 3 (1956).
- <sup>60</sup>M. Challacombe, M. Benzi, and M. M. Tüma (unpublished).
- <sup>61</sup>C. A. White, P. Maslen, M. S. Lee, and M. Head-Gordon, *Chem. Phys. Lett.* **276**, 133 (1997).
- <sup>62</sup>M. Benzi and M. Tüma, Technical Report No. LA-UR-98-2175, Los Alamos National Laboratory, Los Alamos, NM (unpublished).
- <sup>63</sup>M. Benzi, J. Marin, and M. Tüma, in *Fourth IMACS International Symposium on Iterative Methods in Scientific Computation* (IMACS, Austin, TX, 1998).
- <sup>64</sup>F. Liu, *J. Mol. Struct.: THEOCHEM* **230**, 47 (1991).
- <sup>65</sup>L. Cohen and C. Frishberg, *J. Chem. Phys.* **65**, 4234 (1976).
- <sup>66</sup>C. Frishberg, L. Cohen, and P. Blumenau, *Int. J. Quantum Chem., Symp.* **14**, 161 (1980).
- <sup>67</sup>L. Massa and L. Cohen, *Int. J. Quantum Chem., Symp.* **14**, 167 (1980).
- <sup>68</sup>A. Redondo, *Phys. Rev. A* **39**, 4366 (1989).
- <sup>69</sup>A. Redondo and J. C. Marshall, *J. Chem. Phys.* **91**, 5492 (1989).
- <sup>70</sup>C. Waggoner and L. L. Combs, *J. Optim. Theory Appl.* **76**, 225 (1993).
- <sup>71</sup>P. Fantucci and S. Polenzzo, *Int. J. Quantum Chem.* **52**, 817 (1994).
- <sup>72</sup>E. Anderson *et al.*, *LAPACK Users' Guide*, release 2.0 ed. (SIAM, 1994).
- <sup>73</sup>W. H. Press, S. A. Teukolsky, W. T. Vetterling, and B. P. Flannery, *Numerical Recipes in FORTRAN* (Cambridge University Press, Port Chester, NY, 1992).
- <sup>74</sup>L. C. W. Dixon, in *Nonlinear Optimization Theory and Algorithms*, edited by L. C. W. Dixon, E. Spedicato, and G. P. Szegö (Birkhäuser, Boston, 1980), pp. 124–135.
- <sup>75</sup>J. R. Shewchuk, Technical Report No. CMU-CS-94-125, Carnegie Mellon University, Pittsburgh, PA (unpublished).
- <sup>76</sup>L. C. W. Dixon, P. G. Ducksbury, and P. Singh, *J. Optim. Theory Appl.* **47**, 285 (1985).
- <sup>77</sup>J. F. Annett, *Comput. Mater. Sci.* **4**, 23 (1995).
- <sup>78</sup>F. Liu, *J. Mol. Struct.: THEOCHEM* **226**, 197 (1991).
- <sup>79</sup>A. G. Buckley, *Math. Program.* **15**, 200 (1978).
- <sup>80</sup>P. Pulay, *Chem. Phys. Lett.* **73**, 393 (1980).
- <sup>81</sup>R. T. McLay, S. Swift, and G. F. Carey, *J. Par. Dist. Comp.* **37**, 146 (1996).
- <sup>82</sup>S. Goedecker *et al.*, Technical Report No. LA-UR-97-1504, Los Alamos National Laboratory, Los Alamos, NM (unpublished).
- <sup>83</sup>M. R. Field, *SIAM J. Sci. Comput.* **19**, 27 (1998).
- <sup>84</sup>S. Pissanetzky, *Sparse Matrix Technology* (Academic, London, 1984).
- <sup>85</sup>Y. Saad, *Iterative Methods for Sparse Linear Systems* (PWS, Boston, MA, 1996).
- <sup>86</sup>I. S. Duff, A. M. Erisman, and J. K. Reid, *Direct Methods for Sparse Matrices* (Oxford University Press, London, 1986).
- <sup>87</sup>F. G. Gustavson, *ACM Trans. Math. Softw.* **4**, 250 (1978).
- <sup>88</sup>G. Golub and C. F. van Loan, *Matrix Computations* (Johns Hopkins University Press, Baltimore, MD, 1996).
- <sup>89</sup>A. D. Daniels, J. M. Millam, and G. E. Scuseria, *J. Chem. Phys.* **107**, 425 (1997).
- <sup>90</sup>M. Challacombe and E. Schwegler (unpublished).
- <sup>91</sup>P. Pulay, *J. Comput. Chem.* **3**, 556 (1982).
- <sup>92</sup>E. Polak, *Computational Methods in Optimization: A Unified Approach* (Academic, London, 1971).
- <sup>93</sup>M. J. D. Powell, *Math. Program.* **12**, 241 (1977).
- <sup>94</sup>T. R. Gibbs and A. Popolato, *LASL Explosive Property Data* (University of California Press, Berkeley, CA, 1980), p. 109.
- <sup>95</sup>M. J. Frisch *et al.*, GAUSSIAN 94, Revision E.2, Gaussian Inc., Pittsburgh, PA, 1995.
- <sup>96</sup>M. Smotherman, program mm.c, available by anonymous ftp from ftp.nosc.mil under/pub/aburto/mm. Compiled using cc -DN=500 -DUNIX -02 -64 -r10000 and executed as mm -w 50, 1997.
- <sup>97</sup>M. J. Daydé and I. S. Duff, *Lect. Notes Comput. Sci.* **1215**, 108 (1997).